

Sparse Non-Negative Stencils for Anisotropic Diffusion *

Jérôme Fehrenbach[†]

Jean-Marie Mirebeau[‡]

May 22, 2013

Abstract

We introduce a new discretization scheme for Anisotropic Diffusion, AD-LBR, on two and three dimensional cartesian grids. The main features of this scheme is that it is non-negative and has sparse stencils, of cardinality bounded by 6 in 2D, by 12 in 3D, despite allowing diffusion tensors of arbitrary anisotropy. The radius of these stencils is not a-priori bounded however, and can be quite large for pronounced anisotropies. Our scheme also has good spectral properties, which permits larger time steps and avoids e.g. chessboard artifacts.

AD-LBR relies on Lattice Basis Reduction, a tool from discrete mathematics which has recently shown its relevance for the discretization on grids of strongly anisotropic Partial Differential Equations [14]. We prove that AD-LBR is in 2D asymptotically equivalent to a finite element discretization on an anisotropic Delaunay triangulation, a procedure more involved and computationally expensive. Our scheme thus benefits from the theoretical guarantees of this procedure, for a fraction of its cost. Numerical experiments in 2D and 3D illustrate our results.

keywords : Anisotropic Diffusion, Non-Negative Numerical Scheme, Lattice Basis Reduction.

We consider throughout this paper a bounded smooth domain $\Omega \subset \mathbb{R}^d$, where $d \in \{2, 3\}$ denotes the dimension, equipped with a continuous diffusion tensor \mathbf{D} . We do not impose any bound on the diffusion tensor anisotropy, and we are in fact interested in pronounced, non axis-aligned anisotropies. Anisotropic diffusion is here understood in the sense of [25]: the diffusion tensor $\mathbf{D}(z)$, at a point $z \in \Omega$, is a symmetric positive definite matrix whose eigenvalues may have different orders of magnitude. Our results are not relevant for isotropic diffusion with a variable scalar coefficient, as in the pioneering work of Perona and Malik [20].

We address the discretization of the following energy \mathcal{E} , defined for $u \in H^1(\Omega)$:

$$\mathcal{E}(u) := \int_{\Omega} \|\nabla u(z)\|_{\mathbf{D}(z)}^2 dz. \quad (1)$$

We denote $\|e\|_M := \sqrt{\langle e, Me \rangle}$, for any $e \in \mathbb{R}^d$, and any M in the set S_d^+ of symmetric positive definite $d \times d$ matrices. Gradient descent for the energy (1) has the form of a parabolic PDE:

$$\partial_t u = \operatorname{div}(\mathbf{D} \nabla u). \quad (2)$$

This equation, Anisotropic Diffusion, is with its variants at the foundation of powerful image processing techniques. Some variants include curvature terms [19], or diffusion-reaction terms [5]. Time varying and solution dependent diffusion tensors can also be considered. A general exposition can be found in [25], where various choices for the definition of the diffusion tensor \mathbf{D} from the image u , adapted to various applications, are proposed and discussed.

Our contribution in the discretization of the energy (1) results in improved numerical solutions of (2), in terms of accuracy and stability, for a minor increase in complexity. This extends to applications, such as Coherence Enhancing Diffusion and Edge Preserving Diffusion [25], see the numerical experiments in §4, which involve solving (2) using a solution dependent diffusion tensor $\mathbf{D} = \mathbf{D}(u)$. For that purpose, one fixes a time step ΔT , and solves for each integer $n \geq 0$ the linear diffusion equation $\partial_t u = \operatorname{div}(\mathbf{D}_n \nabla u)$ on the interval $[n\Delta T, (n+1)\Delta T]$, with $\mathbf{D}_n := \mathbf{D}(u(n\Delta T))$. In these applications, the diffusion tensor $\mathbf{D}(u)$ is typically defined in terms of the structure tensor [25] of u , in such way that diffusion is pronounced within image homogeneous regions, and *tangentially* along image edges, but not across edges.

In two dimensions, AD-LBR strictly speaking is not the first non-negative scheme for anisotropic diffusion: the proof of Theorem 6 in [25] implicitly defines an alternative 6-point non-negative scheme. This alternative scheme does however lack many of the qualities of AD-LBR: it leads to axis aligned artifacts, spectral aberrations, stencils of larger radius, reduced numerical accuracy, and does not extend to 3D. A detailed description and comparison is presented in §4.1.

*This work was partly supported by ANR grant MESANGE ANR-08-BLAN-0198.

[†]Institut de Mathématiques de Toulouse, Université Paul Sabatier, 31062 TOULOUSE CEDEX 9, France

[‡]CNRS, Laboratory CEREMADE, UMR 7534, University Paris Dauphine, Place du Maréchal De Lattre De Tassigny 75775 PARIS CEDEX 16, France

Consider a scale parameter $h > 0$, and a sampling Ω_h of the domain Ω on the cartesian grid \mathbb{Z}^d , rescaled by h : with obvious notations

$$\Omega_h := \Omega \cap h\mathbb{Z}^d.$$

We introduce a novel discretization of the energy (1), referred to as AD-LBR (Anisotropic Diffusion using Lattice Basis Reduction). It is a sum of squared differences of a discrete map $u \in L^2(\Omega_h)$

$$\mathcal{E}_h(u) := h^{d-2} \sum_{z \in \Omega_h} \sum_{e \in V(z)} \gamma_z(e) |u(z+he) - u(z)|^2 \quad (3)$$

The stencils $V(z) \subset \mathbb{Z}^d$, $z \in \Omega_h$, are symmetric and have cardinality at most 6 in 2D, 12 in 3D. The coefficients $\gamma_z(e) \geq 0$ are non-negative. They are constructed using a classical tool from discrete mathematics, Lattice Basis Reduction, which allows to cheaply build efficient stencils for grid discretizations of Partial Differential Equations (PDEs) involving strongly anisotropic diffusion tensors or Riemannian metrics. This approach has been applied to anisotropic static Hamilton-Jacobi PDEs in [14], resulting in a new numerical scheme: Fast Marching using Lattice Basis Reduction (FM-LBR). Substantial improvements were obtained in comparison with earlier methods, in terms of both accuracy and complexity.

The paper is organized as follows. We describe the stencils of the two dimensional AD-LBR in §1, and state our main 2D result: the asymptotic equivalence of AD-LBR with a finite element discretization on an Anisotropic Delaunay Triangulation. Section §2 provides additional details on the two dimensional stencils of AD-LBR, and describes the three dimensional ones. The more technical §3 details the proof of the 2D equivalence result stated in §1. Two and three dimensional numerical experiments are presented in §4, including qualitative and quantitative comparisons with five other numerical schemes.

1 Description of the scheme, and main results

Our numerical scheme, Anisotropic Diffusion using Lattice Basis Reduction (AD-LBR), involves the construction of stencils whose geometry is tailored after the local diffusion tensor. Its essential feature is non-negativity: the discrete energy $\mathcal{E}_h(u)$ is written as a sum (3) of squared differences of values of u , with non-negative weights $\gamma_z(e) \geq 0$. This discretization is consistent if for each $z \in \Omega_h$, and any smooth u ,

$$h^d \|\nabla u(z)\|_{\mathbf{D}(z)}^2 = h^{d-2} \sum_{e \in V(z)} \gamma_z(e) \langle \nabla u(z), he \rangle^2. \quad (4)$$

Indeed, the left hand side approximates the contribution of the “voxel” $z + [-h/2, h/2]^d$ to the integral (1), while

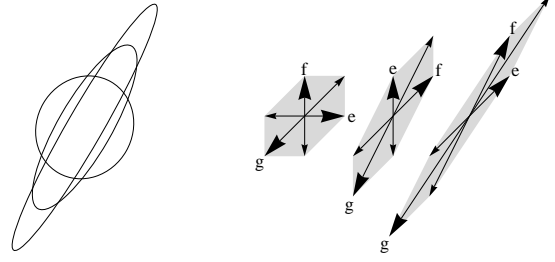


Figure 1: Right: the stencils associated to three matrices M of anisotropy ratios $\kappa(M)$ equal to 1.1, 3.5, 8 respectively. The ellipses $\{z \in \mathbb{R}^2; \|z\|_M = 1\}$ are shown left; their principal axis is aligned with $(\cos(\pi/3), \sin(\pi/3))$. More stencils are shown in [14].

the right hand side is obtained by inserting the first order approximation $u(z+he) \simeq u(z) + \langle \nabla u(z), he \rangle$ in (3). The identity (4) is in turn equivalent to

$$\mathbf{D}(z) = \sum_{e \in V(z)} \gamma_z(e) ee^T. \quad (5)$$

The next lemma shows how to obtain such a decomposition in 2D. We denote by $u^\perp := (-b, a)$ the rotation of a vector $u = (a, b) \in \mathbb{R}^2$ by $\pi/2$, in such way that for all $v \in \mathbb{R}^2$:

$$\langle u^\perp, v \rangle = \det(u, v).$$

Lemma 1. *Let $e_0, e_1, e_2 \in \mathbb{R}^2$ be such that $e_0 + e_1 + e_2 = 0$, and $|\det(e_1, e_2)| = 1$. Then for any $D \in S_2^+$, with the convention $e_{3+i} := e_i$:*

$$D = - \sum_{0 \leq i \leq 2} \langle e_{i+1}^\perp, De_{i+2}^\perp \rangle e_i e_i^T. \quad (6)$$

Proof. Note that $1 = |\det(e_2, e_0)| = |\det(e_0, e_1)|$. Denoting by D' the right hand side of (6), we obtain

$$\begin{aligned} \langle e_1^\perp, D'e_1^\perp \rangle &= -\langle e_0^\perp, De_1^\perp \rangle \langle e_2, e_1^\perp \rangle^2 - \langle e_1^\perp, De_2^\perp \rangle \langle e_0, e_1^\perp \rangle^2 \\ &= -\langle e_0^\perp + e_2^\perp, De_1^\perp \rangle = \langle e_1^\perp, De_1^\perp \rangle. \end{aligned}$$

Thus $\|e_1^\perp\|_{D'} = \|e_1^\perp\|_D$. Likewise $\|e_2^\perp\|_{D'} = \|e_2^\perp\|_D$, and $\|e_1^\perp + e_2^\perp\|_{D'} = \|e_0^\perp\|_{D'} = \|e_0^\perp\|_D = \|e_1^\perp + e_2^\perp\|_D$. Since (e_1^\perp, e_2^\perp) is a basis of \mathbb{R}^2 , the result follows. \square

The diffusion tensor \mathbf{D} is meant to measure gradients, as in (1). In order to measure angles between vectors, we introduce a Riemannian metric $^1 \mathbf{M}$ on the domain Ω , which is proportional to the inverse of \mathbf{D} : for all $z \in \Omega$

$$\mathbf{M}(z) := \mathbf{d}(z) \mathbf{D}(z)^{-1}, \quad \text{where } \mathbf{d}(z) := \det(\mathbf{D}(z))^{1/d}. \quad (7)$$

¹The Laplace Beltrami operator associated to \mathbf{M} does *not* coincide with $\text{div}(\mathbf{D} \nabla \cdot)$, unless \mathbf{D} is identically of determinant 1. This is not an issue for our application.

The normalizing factor $\mathbf{d}(z)$ was chosen so as to normalize the metric determinant: $\det(\mathbf{M}(z)) = 1$. This normalization reflects the fact that the construction of our stencil $V(z)$ depends on the preferred direction of diffusion, and on the amount of anisotropy, whereas the absolute strength of diffusion is irrelevant. In dimension $d = 2$, one easily checks that for any $z \in \Omega$ and any $e, f \in \mathbb{R}^2$, one has

$$\langle e^\perp, \mathbf{D}(z)f^\perp \rangle = \mathbf{d}(z)\langle e, \mathbf{M}(z)f \rangle. \quad (8)$$

The AD-LBR is based on decompositions (5), given by the previous lemma, with a family of vectors $(e_i)_{i=0}^2$ chosen so that the scalar products appearing in (6) are non-positive. The adequate concept is that of M -obtuse superbase of \mathbb{Z}^d [4].

Definition 1. • A basis of \mathbb{Z}^d is a family $(e_i)_{i=1}^d$ of elements of \mathbb{Z}^d such that $|\det(e_1, \dots, e_d)| = 1$.

- A superbase of \mathbb{Z}^d is a family $(e_i)_{i=0}^d$ such that $e_0 + \dots + e_d = 0$, and $(e_i)_{i=1}^d$ is a basis of \mathbb{Z}^d .

Definition 2. Let $M \in S_d^+$. A family $(e_i)_{i \in I}$ of vectors in \mathbb{R}^d is said to be M -obtuse if $\langle e_i, Me_j \rangle \leq 0$ for all distinct $i, j \in I$.

In dimension $d \leq 3$, there exists for each $M \in S_d^+$ at least one M -obtuse superbase of \mathbb{Z}^d [4]. The practical construction of such superbases is discussed in §2, and based on lattice basis reduction algorithms described in [11, 23, 17] (hence the name of our numerical scheme). This construction has a logarithmic numerical cost $\mathcal{O}(\ln \kappa(M))$ in the anisotropy ratio of the matrix M :

$$\kappa(M) := \max_{|u|=|v|=1} \frac{\|u\|_M}{\|v\|_M} = \sqrt{\|M\| \|M^{-1}\|}. \quad (9)$$

The AD-LBR energy $\mathcal{E}_h : L^2(\Omega_h) \rightarrow \mathbb{R}_+$, see (3), is in two dimensions written in terms of the following stencils and coefficients. Let $z \in \Omega$, and let e_0, e_1, e_2 be an $\mathbf{M}(z)$ -obtuse superbase of \mathbb{Z}^2 . We set

$$V(z) := \{e_0, e_1, e_2, -e_0, -e_1, -e_2\}, \quad (10)$$

and for $0 \leq i \leq 2$, with the convention $e_{i+3} := e_i$,

$$\gamma_z(\pm e_i) := -\frac{1}{2} \langle e_{i+1}^\perp, \mathbf{D}(z)e_{i+2}^\perp \rangle. \quad (11)$$

Lemma 1 implies the announced decomposition (5), and the weights γ_z are non-negative in view of (8). These weights $\gamma_z : \mathbb{Z}^2 \rightarrow \mathbb{R}_+$, extended by 0 outside $V(z)$, do not depend on the choice of $\mathbf{M}(z)$ -obtuse superbase (e_0, e_1, e_2) , see Lemma 11. Stencils of the three dimensional AD-LBR are described in §2, and involve a construction of Selling²

²The authors would like to thank Professor P. Q. Nguyen for pointing out this 12 points 3D stencil, which is simpler and sparser than the 14 points stencil proposed by the authors in an earlier version of the manuscript.

[22]. The above description of the stencils $V(z)$ is suitable for periodic, reflected, and Dirichlet boundary conditions (extending u by zero outside Ω_h in the latter case). In the case of Neumann boundary conditions, a slight modification is in order:

$$V(z; h) := \{e \in V(z); z + he \in \Omega_h\}.$$

We have so far established three strongpoints of the AD-LBR:

Non-negativity. Off diagonal coefficients of the symmetric semi-definite $N \times N$ matrix, $N = \#(\Omega_h)$, associated to the energy \mathcal{E}_h are non-positive, while diagonal coefficients are positive.

Sparsity. Stencil cardinality is uniformly bounded, without restriction on the anisotropy ratio $\kappa(\mathbf{D}(z))$ of the diffusion tensor.

Complexity. The construction of the stencil $V(z)$, and of the associated coefficients γ_z , has a logarithmic cost $\mathcal{O}(\ln \kappa(\mathbf{D}(z)))$ in the anisotropy ratio of the diffusion tensor.

The next result, Theorem 1, restricted to the two dimensional case, establishes that AD-LBR is asymptotically equivalent to a more involved and computationally intensive procedure: a finite element discretization of the energy (1), on an *Anisotropic Delaunay Triangulation* (ADT, see [10] and below) of the domain Ω . Under the assumptions of Theorem 1, AD-LBR benefits from two additional guarantees, that we state informally and without proof.

No chessboard artifacts. Some numerical schemes for anisotropic diffusion suffer from chessboard artifacts, in the sense that periodic artifacts develop at the pixel level. Such artifacts cannot develop in finite element discretizations, since they would lead to high frequency oscillations of the finite element interpolant, and therefore to an increase of the energy (12). The asymptotic equivalence of the AD-LBR with a finite element discretization also rules out these defects.

Spectral correctness. The n -th smallest eigenvalue $\lambda_n(h)$ of the symmetric matrix associated to $h^{-d}\mathcal{E}_h$ (3), converges as $h \rightarrow 0$ towards the n -th smallest eigenvalue λ_n of the continuous operator $-\operatorname{div}(\mathbf{D}\nabla)$, for any given integer $n \geq 0$. See Figure 8 page 15 for an illustration. This follows from a similar property of the finite element energy \mathcal{E}'_h (12), and from the asymptotic equivalence (13).

Our convergence result, Theorem 1 below, is specialized to the case of a square periodic domain, which covers reflecting boundary conditions frequently used in image

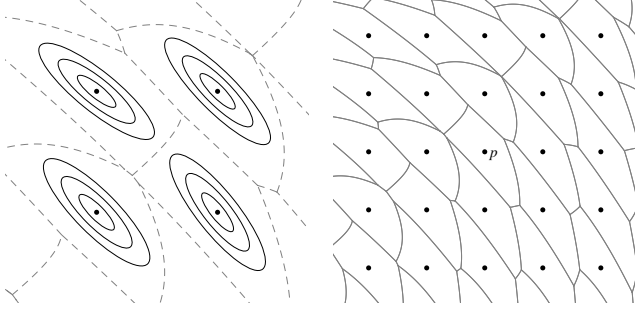


Figure 2: The distance $\delta_p(q)$, from a grid point p to $q \in \mathbb{R}^2$, is defined in terms of the local metric $\mathbf{M}(p)$, see (14). The level lines $\{q \in \mathbb{R}^2; \delta_p(q) = r\}$ are ellipses (left). The collection of points $q \in \mathbb{R}^2$ closer to p than to any other grid point is the Voronoi region of p (left: the boundaries of Voronoi regions are shown dashed). The Voronoi diagram (right) is the collection of all Voronoi regions.

processing. Since the grid discretization must be compatible with the boundary conditions, any scale parameter h appearing in the rest of the paper is assumed to be the inverse of a positive integer:

$$h \in \{1/n; n \geq 1\}.$$

Theorem 1. *Let Ω be the unit square $[0, 1]^2$, equipped with periodic boundary conditions. Let $\mathbf{D} : \Omega \rightarrow S_2^+$ be a (periodic) diffusion tensor with Lipschitz regularity, and let \mathbf{M} be the Riemannian metric defined by (7). When h is sufficiently small, the periodic Riemannian domain (Ω, \mathbf{M}) admits an Anisotropic Delaunay Triangulation \mathcal{T}_h , with collection of vertices $\Omega_h := \Omega \cap h\mathbb{Z}^2$. For $u \in L^2(\Omega_h)$, define*

$$\mathcal{E}'_h(u) := \int_{\Omega} \|\nabla(I_{\mathcal{T}_h} u)(z)\|_{\mathbf{D}(z)}^2 dz, \quad (12)$$

where $I_{\mathcal{T}}$ denotes the piecewise linear interpolation operator on a triangulation \mathcal{T} . Then for some constant $c = c(\mathbf{D})$, independent of u and h ,

$$(1 - ch)\mathcal{E}_h(u) \leq \mathcal{E}'_h(u) \leq (1 + ch)\mathcal{E}_h(u). \quad (13)$$

Let us mention that the finite element discretization on an ADT is a more general procedure than AD-LBR, since it does not require the domain Ω to be sampled on a grid. This flexibility can be used to locally increase the density of vertices, in places where solution u is expected to be less regular, or to insert vertices exactly on $\partial\Omega$ for a better discretization of boundary conditions. (Such refinements are however generally incompatible with image processing since the unknowns, the pixel values, lie by construction on a fixed and given cartesian grid.) Here and as often, the performance of AD-LBR is at the cost of its specialization.

The proof of Theorem 1 is postponed to §3, but for the sake of concreteness, we describe here the concept of

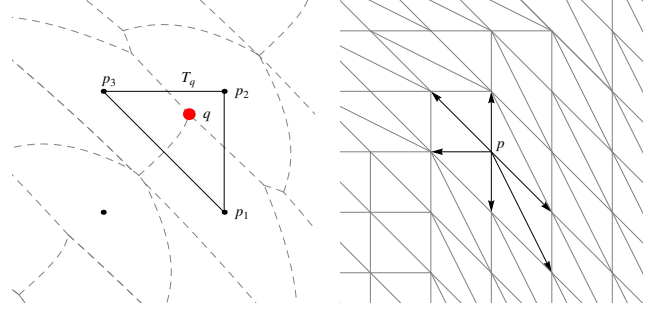


Figure 3: A Voronoi vertex q is a point where the Voronoi regions of at least three grid points $(p_i)_{i=1}^k$ intersect (left: Voronoi region boundaries are shown dashed); here $k = 3$ (values $k > 3$ are non-generic). The dual Voronoi cell T_q , generically a triangle, is the convex envelope of the grid points $\{p_i; 1 \leq i \leq k\}$ (left). The collection of dual Voronoi cells T_q defines a polygonization \mathcal{Q}_h , generically a triangulation, and the Anisotropic Delaunay Triangulation \mathcal{T}_h is obtained by arbitrarily triangulating (if necessary) the elements of \mathcal{Q}_h . Here $\mathcal{T}_h = \mathcal{Q}_h$ (right). The stencil $V_h(p)$, of a vertex p of \mathcal{T}_h , see (17), is represented by arrows (right).

Anisotropic Delaunay Triangulation (ADT) [10]. In the rest of this introduction, and in §3, we assume as in Theorem 1 that the diffusion tensor \mathbf{D} is defined on the square $[0, 1]^2$ and satisfies periodic boundary conditions. We extend it, as well as the metric \mathbf{M} , to the whole plane \mathbb{R}^2 by periodicity.

We specialize the concept of ADT [10], to the domain \mathbb{R}^2 and the collection of vertices $h\mathbb{Z}^2$. For that purpose, we introduce some notations. For all $p, q \in \mathbb{R}^2$, we denote by $\delta_p(q)$ the distance from p to q , as measured by the metric at the point p :

$$\delta_p(q) := \|q - p\|_{\mathbf{M}(p)}. \quad (14)$$

We denote by $\Delta_h(q)$ the least distance from a point $q \in \mathbb{R}^2$, to the grid $h\mathbb{Z}^2$:

$$\Delta_h(q) := \min_{p \in h\mathbb{Z}^2} \delta_p(q). \quad (15)$$

We introduce the Voronoi cell $\text{Vor}_h(p)$ of a grid point $p \in h\mathbb{Z}^2$, which is the collection of points $q \in \mathbb{R}^2$ closer to p than to any other grid point:

$$\text{Vor}_h(p) := \{q \in \mathbb{R}^2; \delta_p(q) = \Delta_h(q)\}. \quad (16)$$

The collection of Voronoi cells is referred to as the Voronoi diagram, see Figure 2. A Voronoi vertex is a point $q \in \mathbb{R}^2$ at which at least three distinct Voronoi regions intersect: $(\text{Vor}_h(p_i))_{i=1}^k$, $k \geq 3$, $p_i \in h\mathbb{Z}^2$. We attach to q a dual Voronoi cell T_q , defined as the convex hull of the points $(p_i)_{i=1}^k$, see Figure 3.

The geometric dual \mathcal{Q}_h , of the Voronoi diagram, is defined as the collection of all dual Voronoi cells T_q . Note that, generically on the metric \mathbf{M} , no more than three Voronoi regions can intersect at any point in \mathbb{R}^2 , thus the elements of \mathcal{Q}_h are generically triangles. If h is small enough, we show in §3 (using the Dual Triangulation Theorem in [10]) that T_q is a strictly convex polygon, of vertices $(p_i)_{i=1}^k$ with the above notations, and that \mathcal{Q}_h is a polygonization (generically a triangulation) of \mathbb{R}^2 , with vertices $h\mathbb{Z}^2$.

Since the metric \mathbf{M} and the vertices $h\mathbb{Z}^2$ are periodic (recall that $h = 1/n$ for some integer $n \geq 1$), arbitrarily triangulating the elements of \mathcal{Q}_h , respecting periodicity, yields a periodic triangulation \mathcal{T}_h .

Definition 3 (ADT, Labelle and Shewchuk [10]). *The triangulation \mathcal{T}_h obtained by the above construction is referred to as an ADT of the domain \mathbb{R}^2 , with collection of vertices $h\mathbb{Z}^2$, and underlying Riemannian metric \mathbf{M} . Since \mathcal{T}_h is \mathbb{Z}^2 -periodic, we also regard it as an ADT of the periodic unit square Ω .*

We establish in §3.1 the existence of the ADT \mathcal{T}_h . Incidentally, we show in Lemma 7 (iii) page 8 that the angles of the elements of \mathcal{T}_h , measured with respect to the local metric \mathbf{M} , are asymptotically acute. This geometrical property (which holds thanks to our special choice of triangulation vertices, on a grid) is linked to the non-negativity of AD-LBR: indeed, it is known that finite elements discretizations such as (12) yield non-negative numerical schemes, and the discrete maximum principle, if the mesh satisfies a non-obtuse angle condition, see Lemma 3.1 in [9].

Subsection §3.2 is devoted to the study of M -obtuse superbases of \mathbb{Z}^2 , and their cousins M -reduced bases of \mathbb{Z}^2 , on which the AD-LBR relies: we discuss their characterization, uniqueness and stability properties. We study in §3.3 the finite element stencils, defined for $p \in h\mathbb{Z}^2$ by

$$V_h(p) := \{e \in \mathbb{Z}^2; [p, p + he] \text{ is an edge of } \mathcal{T}_h\}, \quad (17)$$

see Figure 3 (right). We show that $V_h(p)$ coincides with the AD-LBR stencil $V(p)$, unless the lattice \mathbb{Z}^2 admits a basis *almost orthogonal* with respect to the scalar product associated to $\mathbf{M}(p)$, see Lemma 13. This is tied to the fact that orthogonal grids admit several (usual) Delaunay triangulations. Overcoming this technical difficulty, we conclude the proof of Theorem 1.

Note that the construction of the ADT \mathcal{T}_h is not easy to parallelize, in particular when anisotropy is pronounced since the Voronoi regions of far away points interact. The construction of \mathcal{T}_h also involves solving polynomial equations of degree four, because Voronoi regions boundaries are conics, and Voronoi vertices must be identified at their intersections. In contrast, the AD-LBR stencils are independent of each other, and the numerical cost of their

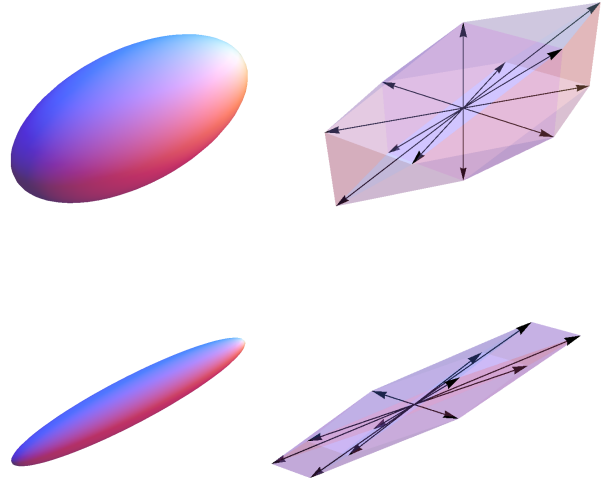


Figure 4: Right: stencil of the AD-LBR, for a symmetric matrix of eigenvector M of anisotropy ratio $\kappa(M)$ equal to 2 (top) or 6 (bottom). The anisotropy is of “needle” type: the two largest eigenvalues of M are equal, and the needle orientation is given by the vector $(4, 2, 3)$. The ellipsoid $\{z \in \mathbb{R}^3; \|z\| \leq 1\}$ is shown left.

construction only grows logarithmically with the metric anisotropy.

2 Construction of obtuse superbases, and three dimensional stencils

Algorithms for the construction of privileged bases of lattices, consisting of short and almost orthogonal vectors, have attracted an important research effort from the mathematical community, over a long period of time. The first such algorithm dates back to Lagrange [11], and is restricted to two dimensional lattices. Methods for high dimensional lattices, such as the LLL algorithm [12], are of key importance for integer programming and cryptography [18]. AD-LBR is based on the original algorithm of Lagrange [11], and on its recent extension to three dimensional lattices [23, 17]. These methods output a basis of \mathbb{Z}^d reduced in the sense of Minkowski, which in dimension $d \leq 4$ is equivalent to the following definition.

We denote by $e_1\mathbb{Z} + \dots + e_k\mathbb{Z}$ the sub-lattice of \mathbb{Z}^d generated by vectors $e_1, \dots, e_k \in \mathbb{Z}^d$. This sub-lattice equals $\{0\}$ by convention if $k = 0$.

Definition 4. *An M -reduced basis of \mathbb{Z}^d , where $d \leq 4$*

and $M \in S_d^+$, is a basis (e_1, \dots, e_d) of \mathbb{Z}^d such that

$$\|e_i\|_M = \min\{\|e\|_M; e \in \mathbb{Z}^d \setminus (e_1\mathbb{Z} + \dots + e_{i-1}\mathbb{Z})\}. \quad (18)$$

For each $d \leq 4$, and each $M \in S_d^+$, there exists at least one M -reduced basis [17]. In contrast, there exists $M \in S_5^+$ for which no basis of \mathbb{Z}^d satisfies (18). The norms of the elements $(e_i)_{i=1}^d$ of an M -reduced basis,

$$\lambda_i(M) := \|e_i\|_M, \quad (19)$$

are called the Minkowski minima, and are independent of the choice of M -reduced basis. In particular, e_1 is the shortest vector of \mathbb{Z}^d , with respect to the norm $\|\cdot\|_M$, and e_2 is the shortest linearly independent vector.

Lemma 2. For any $M \in S_d^+$, $1 \leq i \leq d$,

$$\|M^{-\frac{1}{2}}\|^{-\frac{1}{2}} \leq \lambda_i(M) \leq \|M\|^{\frac{1}{2}}.$$

Proof. Note that $\|M^{-\frac{1}{2}}\|^{-\frac{1}{2}}\|e\| \leq \|e\|_M \leq \|M\|^{\frac{1}{2}}\|e\|$, for any $e \in \mathbb{R}^d$. In addition: (i) any $e \in \mathbb{Z}^d \setminus \{0\}$ satisfies $\|e\| \geq 1$, and (ii) the set $\mathbb{Z}^d \setminus (e_1\mathbb{Z} + \dots + e_{i-1}\mathbb{Z})$ appearing in (18) always contains at least one element e of the canonical basis of \mathbb{R}^d , so that $\|e\| \leq 1$. The announced result easily follows. \square

We emphasize that obtaining an M -reduced basis, i.e. solving the minimization problems (18), is both simple and cheap numerically. In dimension $d = 2$, this is the object of Lagrange's algorithm [11] (later rediscovered by Gauss and often erroneously called Gauss's algorithm, see [17]): initialize (e, f) as the canonical basis of \mathbb{Z}^2 , and

$$\text{Do } (e, f) := (f, e - \text{Round}(\langle e, Mf \rangle / \|f\|_M^2) f), \quad (20)$$

while $\|e\|_M > \|f\|_M$.

This algorithm can be regarded as a two dimensional geometrical generalization of greatest common divisor computation. It can be extended to higher dimension and, in dimension up to four, outputs an M -reduced basis after at most $\mathcal{O}(\ln \kappa(M))$ iterations [17], each consisting of $\mathcal{O}(1)$ operations among reals.

The elements of an M -reduced basis are heuristically never very far from being orthogonal, as illustrated by the following lemma.

Lemma 3. Let $M \in S_d^+$, $d \leq 4$, and let (e_1, \dots, e_d) be an M -reduced basis. Then for any $i, j \in \{1, \dots, d\}$,

$$2|\langle e_i, Me_j \rangle| \leq \|e_i\|_M^2. \quad (21)$$

Proof. Since $\|e_k\|_M$ is an increasing function of $k \in \{1, \dots, d\}$, we may assume that $i < j$. It follows from (18) that $\|e_j\|_M \leq \|e_j + e_i\|_M$, and $\|e_j\|_M \leq \|e_j - e_i\|_M$. Squaring these inequalities, and developing the scalar products, we obtain the announced result. \square

Corollary 1. Let $M \in S_2^+$, and let (e, f) be an M -reduced basis such that $\langle e, Mf \rangle \leq 0$. Then (e, f, g) is an M -obtuse superbase of \mathbb{Z}^2 , with $g := -e - f$. In addition

$$\langle e, Mg \rangle \leq -\|e\|_M^2/2, \quad \langle f, Mg \rangle \leq -\|f\|_M^2/2. \quad (22)$$

Proof. The previous Lemma implies $\langle e, M(e + f) \rangle \geq \|e\|_M^2 - |\langle e, Mf \rangle| \geq \frac{1}{2}\|e\|_M^2$. Likewise $\langle f, M(e + f) \rangle \geq \frac{1}{2}\|f\|_M^2$. The result follows. \square

The practical construction of the two dimensional AD-LBR stencil at a point $z \in \Omega$ amounts to (i) compute an $\mathbf{M}(z)$ -reduced basis (e, f) using Lagrange's algorithm (20), (ii) replace f with $-f$, if necessary, so that $\langle e, Mf \rangle \leq 0$, and (iii) define the stencil $V(z)$ and the weights γ_z in terms of the M -obtuse superbase (e, f, g) of \mathbb{Z}^2 , where $g = -e - f$, as described in (10) and (11).

The rest of this section is devoted to the description of the three dimensional AD-LBR stencils. In contrast with the two dimensional case, the construction of the 3D stencil $V(z)$ at a point $z \in \Omega$, involves a $\mathbf{D}(z)$ -obtuse basis, instead of an $\mathbf{M}(z)$ -obtuse basis.

Proposition 1. Let $D \in S_3^+$, and let (e_1, e_2, e_3) be a D -reduced basis. Let $b_i := \varepsilon_i e_{\sigma(i)}$, for all $1 \leq i \leq 3$, where the signs $\varepsilon_1, \varepsilon_2, \varepsilon_3 \in \{-1, 1\}$, and the permutation σ of $\{1, 2, 3\}$ are chosen so that

$$|\langle b_1, Db_2 \rangle| \leq \min\{-\langle b_1, Db_3 \rangle, -\langle b_2, Db_3 \rangle\}. \quad (23)$$

Then the following is a D -obtuse superbase:

$$\begin{cases} (b_1, b_2, b_3, -b_1 - b_2 - b_3) & \text{if } \langle b_1, Db_2 \rangle \leq 0, \\ (-b_1, b_2, b_1 + b_3, -b_2 - b_3) & \text{otherwise.} \end{cases} \quad (24)$$

Proof. To achieve (23), one can choose σ such that $b'_i := e_{\sigma(i)}$ satisfies $|\langle b'_1, Db'_2 \rangle| \leq |\langle b'_1, Db'_3 \rangle| \leq |\langle b'_2, Db'_3 \rangle|$. Then choose the signs $(\varepsilon_i)_{i=1}^3$ such that $b_i := \varepsilon_i b'_i$ satisfies $\langle b_1, Db_3 \rangle \leq 0$ and $\langle b_2, Db_3 \rangle \leq 0$.

The two families of vectors appearing in (24) are clearly superbases. We thus only need to show that they are D -obtuse; in other words that $\langle e, Df \rangle \leq 0$ for any two distinct elements e, f of these families. Note that for all distinct $i, j \in \{1, 2, 3\}$, using (21),

$$2|\langle b_i, Db_j \rangle| \leq \|b_i\|_D^2.$$

In the case where $\langle b_1, Db_2 \rangle \leq 0$, the pairwise scalar products between b_1, b_2, b_3 are non-positive by construction. In addition

$$\begin{aligned} & 2\langle b_1 + b_2 + b_3, Db_1 \rangle \\ & \geq (\|b_1\|_D^2 - 2|\langle b_1, Db_2 \rangle|) + (\|b_1\|_D^2 - 2|\langle b_1, Db_3 \rangle|) \geq 0. \end{aligned}$$

Likewise $\langle b_1 + b_2 + b_3, Db_i \rangle \geq 0$ for all $i \in \{1, 2, 3\}$, which concludes the proof.

We next turn to the second case, where $\langle b_1, Db_2 \rangle \geq 0$. Enumerating all scalar products we obtain

$$\begin{aligned}\langle b_1, D(b_1 + b_3) \rangle &\geq \|b_1\|_D^2 - |\langle b_1, Db_3 \rangle| \geq 0, \\ \langle b_1, D(-b_2 - b_3) \rangle &= -\langle b_1, Db_2 \rangle - \langle b_1, Db_3 \rangle \geq 0, \\ -\langle b_2, D(b_1 + b_3) \rangle &= -\langle b_2, Db_1 \rangle - \langle b_2, Db_3 \rangle \geq 0, \\ \langle b_2, D(b_2 + b_3) \rangle &\geq \|b_2\|_D^2 - |\langle b_2, Db_3 \rangle| \geq 0,\end{aligned}$$

and finally

$$\begin{aligned}2\langle b_1 + b_3, D(b_2 + b_3) \rangle &\geq 2\langle b_1, Db_2 \rangle \\ + (\|b_3\|^2 - 2|\langle b_1, Db_3 \rangle|) + (\|b_3\|^2 - 2|\langle b_2, Db_3 \rangle|) &\geq 0.\end{aligned}$$

This concludes the proof. \square

In view of the previous Proposition, obtaining a D -obtuse superbase of \mathbb{Z}^3 has numerical cost $\mathcal{O}(\ln \kappa(D))$. Indeed a D -reduced basis needs to be computed in a preliminary step, after what Proposition 1 is applied for a negligible $\mathcal{O}(1)$ cost. An alternative method for the construction of D -obtuse superbases of \mathbb{Z}^3 is presented in [4] and in appendix B of [3], but its numerical complexity is not known to the authors.

The three dimensional AD-LBR is defined by the following stencils and coefficients. Let $z \in \Omega$, let $D := \mathbf{D}(z)$, and let $(e_i)_{i=0}^3$ be a D -obtuse superbase of \mathbb{Z}^3 . We set

$$V(z) := \{e_k \times e_l; k, l \in \{0, 1, 2, 3\}, k \neq l\},$$

and if $\{i, j, k, l\} = \{0, 1, 2, 3\}$, $i \neq j$ and $k \neq l$, then

$$\gamma_z(e_k \times e_l) := -\frac{1}{2}\langle e_i, De_j \rangle.$$

As announced, $\#(V(z)) = 12$, and the weights γ_z are non-negative. The proof of the scheme consistency (5), due to Selling [22], is reproduced in the next lemma for completeness. A generalization, appearing in Appendix B of [3], allows in arbitrary dimension to build a non-negative decomposition of the form (25) from a D -obtuse superbase of \mathbb{Z}^d . However the non existence of such a superbase, for some matrices $D \in S_4^+$, forbids a straightforward extension of AD-LBR to higher dimension.

Lemma 4 (Selling [22]). *Let $(e_i)_{i=0}^3$ be a superbase of \mathbb{Z}^3 . For all i, j, k, l such that $\{i, j, k, l\} = \{0, 1, 2, 3\}$, $i < j$, and $k < l$, let $c_{ij} := e_k \times e_l$. Then, for any $D \in S_3^+$:*

$$D = - \sum_{0 \leq i < j \leq 3} \langle e_i, De_j \rangle c_{ij} c_{ij}^T. \quad (25)$$

Proof. Let i, j, k, l be as in the definition of c_{ij} . Then

$$\langle e_i, c_{ij} \rangle = \langle e_i, e_k \times e_l \rangle = \det(e_i, e_k, e_l) \in \{-1, 1\},$$

since (e_i, e_k, e_l) is a basis of \mathbb{Z}^3 . Also

$$\begin{aligned}\langle e_j, c_{ij} \rangle &= \langle -e_i - e_k - e_l, e_k \times e_l \rangle \\ &= -\langle e_i, e_k \times e_l \rangle = -\langle e_i, c_{ij} \rangle.\end{aligned}$$

In addition, clearly, $\langle e_k, c_{ij} \rangle = \langle e_l, c_{ij} \rangle = 0$. Denoting by D' the right hand side of (25), we obtain as a result

$$\begin{aligned}\langle e_0, D'e_0 \rangle &= -\langle e_0, De_1 \rangle - \langle e_0, De_2 \rangle - \langle e_0, De_3 \rangle \\ &= \langle e_0, D(-e_1 - e_2 - e_3) \rangle = \langle e_0, De_0 \rangle. \\ \langle e_0, D'e_1 \rangle &= -\langle e_0, De_1 \rangle \langle e_0, c_{01} \rangle \langle e_1, c_{01} \rangle = \langle e_0, De_1 \rangle.\end{aligned}$$

Likewise $\langle e_i, D'e_j \rangle = \langle e_i, De_j \rangle$ for all $i, j \in \{1, 2, 3, 4\}$. It follows as announced that $D = D'$. \square

3 Equivalence to a finite element discretization

This section is devoted to the proof of Theorem 1: the asymptotic equivalence of AD-LBR with a finite element discretization on an Anisotropic Delaunay Triangulation (ADT). We use the notations of §1. The existence of the ADT \mathcal{T}_h is established in the first subsection, for h sufficiently small, as well as a few of its properties. The second subsection is devoted to the study of M -reduced bases. Theorem 1 is proved in the third subsection, by comparing the stencils of the AD-LBR and of the finite element discretization.

We denote by κ the maximum anisotropy ratio (9) of the diffusion tensor

$$\kappa := \max_{z \in \Omega} \kappa(\mathbf{D}(z)). \quad (26)$$

Observing that $\kappa(\mathbf{D}(z)) = \kappa(\mathbf{M}(z))$, and recalling that $\det(\mathbf{M}(z)) = 1$, one easily checks that

$$\kappa^{-\frac{1}{2}} \|e\| \leq \|e\|_{\mathbf{M}(z)} \leq \kappa^{\frac{1}{2}} \|e\|, \quad (27)$$

for all $z \in \Omega$ and all $e \in \mathbb{R}^2$.

3.1 Existence of an ADT

Our first lemma provides a uniform bound on the size of the Voronoi regions, see Figure 3, involved in the construction of the ADT.

Lemma 5. (i) *For all $r \in \mathbb{R}^2$, one has $\Delta_h(r) \leq \kappa^{\frac{1}{2}} h$.*

(ii) *If $p, q \in h\mathbb{Z}^2$, and $r \in \text{Vor}_h(p) \cap \text{Vor}_h(q)$, then $\|p - r\| \leq \kappa h$ and $\|p - q\| \leq 2\kappa h$.*

Proof. Point (i). Rounding the coordinates of r to a nearest multiple of h , we obtain a point $p \in h\mathbb{Z}^2$ such that $\|p - r\| \leq h$. Recalling (27) we obtain $\delta_p(r) \leq \kappa^{\frac{1}{2}} h$, and therefore $\Delta_h(r) \leq \kappa^{\frac{1}{2}} h$ in view of (15).

Point (ii). We have $\kappa^{-\frac{1}{2}} \|p - r\| \leq \delta_p(r) = \Delta_h(r) \leq \kappa^{\frac{1}{2}} h$. Thus $\|p - r\| \leq \kappa h$, and likewise $\|q - r\| \leq \kappa h$. Finally, by the triangle inequality, $\|p - q\| \leq \|p - r\| + \|q - r\| \leq 2\kappa h$. \square

Following the notations of [10], we denote by $\tau(p, q)$, $p, q \in \mathbb{R}^2$, the smallest constant $\tau \geq 1$ such that

$$\tau^{-1} \delta_p(r) \leq \delta_q(r) \leq \tau \delta_p(r), \quad \text{for all } r \in \mathbb{R}^2.$$

Equivalently, in the sense of symmetric matrices,

$$\tau^{-2} \mathbf{M}(p) \leq \mathbf{M}(q) \leq \tau^2 \mathbf{M}(p). \quad (28)$$

We also define a quantity $\tau_h \geq 1$, closely related to the modulus of continuity of the metric \mathbf{M} :

$$\tau_h := \max\{\tau(p, q); \|p - q\| \leq 2\kappa h\}. \quad (29)$$

One has $\tau_h \rightarrow 1$ as $h \rightarrow 0$, for any continuous metric \mathbf{M} (indeed \mathbf{M} is periodic and therefore uniformly continuous). If \mathbf{M} is Lipschitz, as assumed in Theorem 1, then $\tau_h = 1 + \mathcal{O}(h)$.

We show in the next lemma the existence of an ADT, by applying the main result of [10], under the assumption that τ_h is sufficiently small. More precisely, we assume in the rest of this subsection that

$$\tau_h < \sqrt{1 + \kappa^{-2}}. \quad (30)$$

Lemma 6. (i) If $p, q \in h\mathbb{Z}^2$, $p \neq q$, and $r \in \text{Vor}_h(p) \cap \text{Vor}_h(q)$, then $\delta_p(r) < \delta_q(r) / \sqrt{\tau(p, q)^2 - 1}$.

(ii) The geometric dual \mathcal{Q}_h of the Voronoi diagram is, as announced in §1, a polygonization of \mathbb{R}^2 into strictly convex polygons, with vertices $h\mathbb{Z}^2$.

Proof. Point (i). We may assume that $\tau(p, q) > 1$, otherwise there is nothing to prove. Point (ii) of Lemma 5 implies that $\|p - q\| \leq 2\kappa h$, thus

$$\sqrt{\tau(p, q)^2 - 1} \leq \sqrt{\tau_h^2 - 1} < \kappa^{-1}.$$

On the other hand $\delta_p(q) \geq \kappa^{-\frac{1}{2}} \|q - p\| \geq \kappa^{-\frac{1}{2}} h$, and $\delta_p(r) \leq \Delta_h(r) \leq \kappa^{\frac{1}{2}} h$. The announced inequality follows.

Point (ii). We apply Theorem 7 (Dual Triangulation Theorem) in [10]. Since the domain \mathbb{R}^2 has no boundary, it suffices to check that all the Voronoi arcs and vertices are *wedged*, see [10]. This condition means that for any $p, q \in h\mathbb{Z}^2$ such that $p \neq q$, and any $r \in \text{Vor}_h(p) \cap \text{Vor}_h(q)$, one has $(r - q) \mathbf{M}(q)(p - q) > 0$, and likewise exchanging the roles of p and q . Heuristically, it expresses the acuteness of some angles measured in the local metric. Lemma 5 in [10] shows that this condition follows from point (i) of this lemma, which concludes the proof. \square

We recall that \mathcal{T}_h is the triangulation obtained by arbitrarily triangulating the polygonization \mathcal{Q}_h of the previous lemma, respecting periodicity, see Definition 3. Generically \mathcal{Q}_h is already a triangulation, hence $\mathcal{T}_h = \mathcal{Q}_h$, see §1. The Voronoi regions Vor_h , and the triangulation \mathcal{T}_h , are illustrated in Figures 2 and 3.

The next lemma provides estimates of the diameter, the area, and the angles of the elements of \mathcal{T}_h . These geometrical properties also have an interpretation in the context of lattices: (ii) shows that the edges of any triangle $T \in \mathcal{T}_h$ define a superbase (e, f, g) of \mathbb{Z}^2 , and (iii) that this superbase is *almost* $\mathbf{M}(z)$ -obtuse, for any $z \in T$.

Note that the vertices p, q, r of any triangle $T \in \mathcal{T}_h$ satisfy by construction

$$\text{Vor}_h(p) \cap \text{Vor}_h(q) \cap \text{Vor}_h(r) \neq \emptyset. \quad (31)$$

Lemma 7. Denote by he, hf, hg the edges of a triangle $T \in \mathcal{T}_h$, where $e, f, g \in \mathbb{Z}^2$ are oriented so that $e + f + g = 0$. Then

- (i) $\max\{\|e\|, \|f\|, \|g\|\} \leq 2\kappa$.
- (ii) $|\det(e, f)| = 1$, thus $|T| = h^2/2$.
- (iii) $\langle e, \mathbf{M}(z)f \rangle \leq \theta_h$, for any $z \in T$, where $\theta_h \rightarrow 0$ as $h \rightarrow 0$. (Explicitly: $\theta_h = \kappa(3 + 9\tau_{2h}^2)(\tau_{2h}^2 - 1)$)

Proof. Point (i). We denote by p, q, r the vertices of T , ordered in such way that $p + he = q$, $q + hf = r$, $r + hg = p$. The announced estimate follows from (31), and from point (ii) of Lemma 5.

Point (ii). Since \mathcal{T}_h is a conforming triangulation, the intersection of T with the collection $h\mathbb{Z}^2$ of all vertices of \mathcal{T}_h consists of only three points: the vertices p, q, r of T . Thus the triangle of vertices $-e, 0, f$, homothetic to T , contains no point of integer coordinates but its vertices. This implies that (e, f) is a basis of \mathbb{Z}^2 , hence $|\det(e, f)| = 1$, as announced.

Point (iii). The pairwise distances between p, q, r are bounded by $2\kappa h$, see point (i), and since $z \in T$ so are the pairwise distances between p, q, r, z . Defining $s := p - q + r \in h\mathbb{Z}^2$, and observing that $\|s - p\| = \|r - q\| \leq 2\kappa h$, we find that the pairwise distances between p, q, r, z, s are bounded by $4\kappa h$.

Let $x \in \text{Vor}_h(p) \cap \text{Vor}_h(q) \cap \text{Vor}_h(r)$. We have $\delta_p(x) = \delta_q(x) = \delta_r(x) = \Delta_h(x) \leq \delta_s(x)$, thus

$$\delta_s(x)^2 \geq \delta_p(x)^2 - \delta_q(x)^2 + \delta_r(x)^2. \quad (32)$$

(For intuition: in a classical Delaunay triangulation, x would be the circumcenter of T , and (32) would state that s is outside the circumcircle of T .) Denoting $M := \mathbf{M}(z)$, and $\delta := \Delta_h(x)$, we obtain

$$\begin{aligned} |\delta_p(x)^2 - \|x - p\|_M^2| &\leq \delta_p(x)^2(\tau(p, z)^2 - 1) \\ &\leq \delta^2(\tau_{2h}^2 - 1), \end{aligned} \quad (33)$$

using Lemma 5, and likewise for q, r . We also have

$$\begin{aligned} \delta_s(x) &= \|p - q + r - x\|_{\mathbf{M}(s)} \\ &\leq \|p - x\|_{\mathbf{M}(s)} + \|q - x\|_{\mathbf{M}(s)} + \|r - x\|_{\mathbf{M}(s)} \\ &\leq 3\delta\tau_{2h}. \end{aligned}$$

Thus, proceeding as in (33),

$$|\delta_s(x) - \|s - x\|_M|^2 \leq \delta_s(x)^2(\tau_{2h}^2 - 1) \leq 9\delta^2\tau_{2h}^2(\tau_{2h}^2 - 1).$$

Inserting in (32) these estimates of $\delta_\star(x)$, $\star \in \{p, q, r, s\}$, and using the fact that $\delta \leq \kappa^{\frac{1}{2}}h$, see Lemma 5, we obtain after expansion the announced estimate of $\langle e, Mf \rangle$. \square

We next rewrite the finite element energy \mathcal{E}'_h (12) in a form similar to that of the AD-LBR energy \mathcal{E}_h (3). Let $\varphi_p^h : \mathbb{R}^2 \rightarrow \mathbb{R}$ be the piecewise linear function on \mathcal{T}_h such that $\varphi_p^h(p) = 1$, and $\varphi_p^h(q) = 0$ for any vertex $q \in h\mathbb{Z}^2$ distinct from p . This is the classical “hat function” encountered in finite element analysis. For all $p \in h\mathbb{Z}^2$, $e \in \mathbb{Z}^2 \setminus \{0\}$, let

$$\gamma_p^h(e) := -\frac{1}{2} \int_{\mathbb{R}^2} \langle \nabla \varphi_p^h(z), \mathbf{D}(z) \nabla \varphi_{p+he}^h(z) \rangle dz \quad (34)$$

Clearly, $\gamma_p^h(e) = 0$ if $[p, p+he]$ is not an edge of \mathcal{T}_h , in other words if e does not belong to the stencil $V_h(p)$, defined in (17). We express in the next lemma the finite element energy \mathcal{E}'_h (12) in terms of the stencils V_h and of the (potentially negative) weights γ_p^h .

Lemma 8. *For any $u \in L^2(\Omega_h)$, extended by periodicity to $h\mathbb{Z}^2$, one has*

$$\mathcal{E}'_h(u) = \sum_{p \in \Omega_h} \sum_{e \in V_h(p)} \gamma_p^h(e) |u(p+he) - u(p)|^2. \quad (35)$$

Proof. For any triangle $T \in \mathcal{T}_h$, and any $p, q \in h\mathbb{Z}^2$, we denote

$$s_T(p, q) := \int_T \langle \nabla \varphi_p^h(z), \mathbf{D}(z) \nabla \varphi_q^h(z) \rangle dz.$$

Clearly $s_T(p, q) = 0$ if q or p is not a vertex of T . The coefficient $\gamma_p^h(e)$, $e \in \mathbb{Z}^2$, is thus given by the following sum with at most two non-zero terms:

$$\gamma_p^h(e) = -\frac{1}{2} \sum_{T \in \mathcal{T}_h} s_T(p, p+he). \quad (36)$$

Let $p, q, r \in h\mathbb{Z}^2$ be the vertices of a triangle $T \in \mathcal{T}_h$. Since the sum $\varphi_p^h + \varphi_q^h + \varphi_r^h$ is constant on T , equal to 1, it has a null gradient on T , and therefore

$$s_T(p, p) + s_T(p, q) + s_T(p, r) = 0.$$

Using this relation, and the two similar ones obtained by a cyclic permutation of p, q, r , we obtain

$$\begin{aligned} & \int_T \|\nabla(\mathbf{I}_{\mathcal{T}_h} u)(z)\|_{\mathbf{D}(z)}^2 dz \\ &= u(p)^2 s_T(p, p) + u(q)^2 s_T(q, q) + u(r)^2 s_T(r, r) \\ & \quad + 2u(p)u(q) s_T(p, q) + 2u(q)u(r) s_T(q, r) \\ & \quad + 2u(r)u(p) s_T(r, p), \\ &= -s_T(p, q)(u(p) - u(q))^2 - s_T(q, r)(u(q) - u(r))^2 \\ & \quad - s_T(r, p)(u(r) - u(p))^2. \end{aligned}$$

Summing this expression over all $T \in \mathcal{T}_h$, and combining it with (36), we obtain (35), which concludes the proof. \square

Finally, we provide an approximation of the coefficients γ_p^h which will be easily compared with the AD-LBR weights γ_p (11).

Lemma 9. *Consider an edge $[p, p+he]$ of \mathcal{T}_h , shared by the two distinct triangles $T, T' \in \mathcal{T}_h$. Let hf, hg (resp. hf', hg') be the two other vector edges of T (resp. T'), oriented so that $e+f+g=0$ (resp. $e+f'+g'=0$). Then*

$$\left| \gamma_p^h(e) + \frac{1}{4} (\langle f^\perp, \mathbf{D}(p) g^\perp \rangle + \langle f'^\perp, \mathbf{D}(p) g'^\perp \rangle) \right| \leq \varepsilon_h,$$

where $\varepsilon_h := 2\kappa^2 \max\{\|\mathbf{D}(x) - \mathbf{D}(y)\|; \|x - y\| \leq 2\kappa h\}$.

Proof. We assume, up to exchanging f and g , that $[p, p-hg]$ is an edge of T . Let $\alpha := \det(e, f) \in \{-1, 1\}$, see point (ii) of Lemma 7; note that $\alpha = \det(f, g) = \det(g, e)$. Let γ be the constant value of $\nabla \varphi_p^h$ on T . Then $\langle \gamma, he \rangle = -1$ and $\langle \gamma, hg \rangle = 1$. These two independent linear identities are also satisfied by $\alpha f^\perp/h$, hence $\nabla \varphi_p^h = \gamma = \alpha f^\perp/h$ on T .

Denoting $q := p+hg$, we obtain likewise $\nabla \varphi_q^h = \alpha g^\perp/h$ on T . Hence recalling that $|T| = h^2/2$:

$$\begin{aligned} \int_T \langle \nabla \varphi_p^h, \mathbf{D}(p) \nabla \varphi_q^h \rangle &= \frac{h^2}{2} \left\langle \frac{\alpha f^\perp}{h}, \mathbf{D}(p) \frac{\alpha g^\perp}{h} \right\rangle \\ &= \frac{1}{2} \langle f^\perp, \mathbf{D}(p) g^\perp \rangle \end{aligned}$$

Therefore, using point (i) of Lemma 7 in the last step,

$$\begin{aligned} & \left| \int_T \langle \nabla \varphi_p^h, \mathbf{D}(z) \nabla \varphi_q^h \rangle dz - \frac{1}{2} \langle f^\perp, \mathbf{D}(p) g^\perp \rangle \right| \\ &= \left| \int_T \langle \nabla \varphi_p^h, (\mathbf{D}(z) - \mathbf{D}(p)) \nabla \varphi_q^h \rangle dz \right| \\ &\leq \frac{h^2}{2} \frac{2\kappa}{h} \frac{2\kappa}{h} \max\{\|\mathbf{D}(z) - \mathbf{D}(p)\|; z \in T\} \leq \varepsilon_h. \end{aligned}$$

Proceeding likewise on T' , and recalling (34) (or (36)), we conclude the proof. \square

3.2 Some properties of M -reduced bases

We establish some technical properties of M -reduced bases, thanks to which we will be able to compare in §3.3 the “geometric” construction of the ADT finite element stencils V_h , with the lattice based construction of the AD-LBR stencils V .

Lemma 10. *Let $M \in S_2^+$, let $e_1, \dots, e_n \in \mathbb{Z}^2$, $n > 2$, and let $\varepsilon \in \{-1, 1\}$. Assume that for all $1 \leq i \leq n$, with the convention $e_{n+1} := e_1$:*

$$\det(e_i, e_{i+1}) = \varepsilon, \quad (37)$$

$$\langle e_i, M e_{i+1} \rangle > -\frac{1}{2} \min\{\|e_i\|_M^2, \|e_{i+1}\|_M^2\}. \quad (38)$$

Then any M -reduced basis (e, f) satisfies

$$\{e, f\} \subset \{e_1, \dots, e_n\}.$$

Proof. Let $z \in \mathbb{Z}^2 \setminus \{e_1, \dots, e_n\}$. Our objective is to show that z cannot be an element of an M -reduced basis, and we may therefore assume that z has co-prime coordinates.

It follows from (37) that the closed polygonal line of consecutive vertices e_1, \dots, e_n , circles at least once around the origin, see Figure 5. Hence $z = \alpha e_i + \beta e_{i+1}$, for some $1 \leq i \leq n$ and some $\alpha, \beta \geq 0$. Since (e_i, e_{i+1}) is a basis of \mathbb{Z}^2 (indeed $|\det(e_i, e_{i+1})| = 1$), the coefficients α and β are integers. Since $z \notin \{e_i, e_{i+1}\}$, $\alpha + \beta \geq 2$. Since z has co-prime coordinates, $\alpha\beta \neq 0$.

Assuming without loss of generality that $\|e_i\|_M \geq \|e_{i+1}\|_M$, we obtain using (38):

$$\begin{aligned} \|z\|_M^2 &= \alpha^2 \|e_i\|_M^2 + \beta^2 \|e_{i+1}\|_M^2 + 2\alpha\beta \langle e_i, M e_{i+1} \rangle \\ &> \alpha^2 \|e_i\|_M^2 + \beta^2 \|e_{i+1}\|_M^2 - \alpha\beta \min\{\|e_i\|_M^2, \|e_{i+1}\|_M^2\} \\ &\geq \|e_i\|_M^2 + (\alpha^2 + \beta^2 - 1 - \alpha\beta) \|e_{i+1}\|_M^2. \end{aligned}$$

Observing that $\alpha^2 + \beta^2 - 1 - \alpha\beta \geq 0$ for all $\alpha, \beta \in [1, \infty[$, we obtain $\|z\|_M > \|e_i\|_M$. Since e_i and e_{i+1} are linearly independent, we have $\|e_i\|_M \geq \lambda_2(M)$. Finally $\|z\|_M > \lambda_2(M)$, hence z cannot be an element of an M -reduced basis, which concludes the proof. \square

The next corollary reverses the construction, presented in Corollary 1, of an M -obtuse superbase from an M -reduced basis.

Corollary 2. *Let $M \in S_2^+$ and let (e, f, g) be an M -obtuse superbase of \mathbb{Z}^2 , ordered so that $\|e\|_M \leq \|f\|_M \leq \|g\|_M$. Then (e, f) is an M -reduced basis.*

Proof. The family $(e, -g, f, -e, g, -f)$ satisfies by construction the conditions of the previous lemma. Hence any $\mathbf{M}(z)$ -reduced basis (e', f') of \mathbb{Z}^2 satisfies $\{e', f'\} \subset \{e, f, g, -e, -f, -g\}$. Observing that e' and f' are linearly independent, that $\|e'\|_M \leq \|f'\|_M$, and that $\|e\|_M \leq \|f\|_M \leq \|g\|_M$, we obtain that $\|e\|_M \leq \|e'\|_M$ and $\|f\|_M \leq \|f'\|_M$. Recalling that M -reduced bases are defined by the minimality of their $\|\cdot\|_M$ -norms, see Definition 4, we obtain as announced that (e, f) is an M -reduced basis. \square

The previous lemma shows that for any $z \in \Omega$, there exists an $\mathbf{M}(z)$ -reduced basis (e, f) such that

$$V(z) = \{e, f, e + f, -e, -f, -e - f\}. \quad (39)$$

Given $M \in S_2^+$, and an M -reduced basis (e, f) of \mathbb{Z}^2 , we denote $\mu(M) := |\langle e, Mf \rangle|$. This value can be expressed in terms of the Minkowski minima (19) and thus does not depend on the particular choice of M -reduced basis. Indeed, recalling the identity

$$\langle e, Mf \rangle^2 + \det(M) \det(e, f)^2 = \|e\|_M^2 \|f\|_M^2,$$

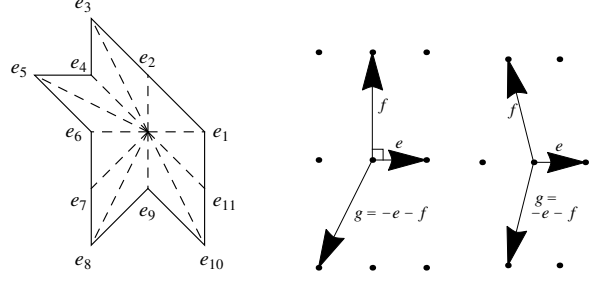


Figure 5: (left) A family e_1, \dots, e_{11} satisfying condition (37) of Lemma 10: the closed polygonal line of vertices (e_1, \dots, e_{11}) circles (at least) once around the origin, and the triangles $(0, e_i, e_{i+1})$ have area $1/2$. (Center and right) The lattice \mathbb{Z}^2 , and an M -reduced basis (e, f) , shown after a linear change of coordinates by A , such that $A^T A = M \in S_2^+$. Case $\mu(M) = 0$ (center), and case $2\mu(M) = \lambda_1(M)^2$ (right).

we obtain

$$\mu(M) = |\langle e, Mf \rangle| = \sqrt{\lambda_1(M)^2 \lambda_2(M)^2 - \det(M)}. \quad (40)$$

In addition one has

$$0 \leq 2\mu(M) \leq \lambda_1(M)^2, \quad (41)$$

where the right hand side follows from Lemma 3. A vanishing value, $\mu(M) = 0$, indicates that the lattice \mathbb{Z}^2 admits an M -orthogonal basis. In contrast, when the upper bound is met, $2\mu(M) = \lambda_1(M)^2$, one has $\|f\|_M = \|f + \varepsilon e\|_M$ for $\varepsilon := -\text{sign}\langle e, Mf \rangle$, hence the reduced basis (e, f) is not unique even up to sign changes. See Figure 5.

We next show that the stencils of the AD-LBR do not depend on the choices of reduced bases, as was announced in the introduction.

Lemma 11. *The weights $\gamma_z : \mathbb{Z}^2 \rightarrow \mathbb{R}_+$ used in the AD-LBR at a point $z \in \Omega$ (defined on $V(z)$ by (11) and extended to \mathbb{Z}^2 by 0), do not depend on the choice of $\mathbf{M}(z)$ -obtuse superbase of \mathbb{Z}^2 .*

Proof. We denote $M := \mathbf{M}(z)$ and $D := \mathbf{D}(z)$. Let (e, f, g) and (e', f', g') be two M -obtuse superbases, and let V, V' and $\gamma, \gamma' : \mathbb{Z}^2 \rightarrow \mathbb{R}_+$ be the corresponding AD-LBR stencils and weights defined by (10) and (11). We may assume, using Corollary 2 and up to reordering, that (e, f) and (e', f') are M -reduced bases.

Corollary 1 states that the scalar products $\langle e, Mg \rangle$, $\langle f, Mg \rangle$, $\langle e', Mg' \rangle$ and $\langle f', Mg' \rangle$ are (strictly) negative. On the other hand

$$\langle e, Mf \rangle = \langle e', Mf' \rangle = -\mu(M) \leq 0. \quad (42)$$

Applying Lemma 10 to the family

$$(e', -g', f', -e', g', -f')$$

we obtain that

$$\{e, f\} \subset \{e', f', g', -e', -f', -g'\}. \quad (43)$$

If $\mu(M) \neq 0$, then $\langle e, Mf \rangle$ and $\langle e', Mf' \rangle$ are negative, and not merely non-positive, thus $\{e, f\} \subset \{e', f', g'\}$, or $\{e, f\} \subset \{-e', -f', -g'\}$. Since $e + f + g = 0 = e' + f' + g'$, it follows that $\{e, f, g\} = \{e', f', g'\}$, or $\{e, f, g\} = \{-e', -f', -g'\}$. The stencils V, V' are thus identical, see (10), and so are the weights γ, γ' .

If $\mu(M) = 0$, then the stencils V, V' may not be identical. Observe however that $\langle e^\perp, Df^\perp \rangle = 0 = \langle e'^\perp, Df'^\perp \rangle$, using (8). Hence using the weights expression (11):

$$\begin{aligned} \gamma(\pm g) &= -\langle e^\perp, Df^\perp \rangle / 2 = 0, \\ \gamma(\pm e) &= \|f^\perp\|_D^2 / 2, \quad \gamma(\pm f) = \|e^\perp\|_D^2 / 2, \end{aligned} \quad (44)$$

and likewise for γ', e', f', g' . Note also that $\|g'\|_M^2 = \|e'\|_M^2 + \|f'\|_M^2 > \lambda_2(M)^2$, hence e and f are different from g' and $-g'$. It follows from (43) that $\{e, f\} = \{\varepsilon_1 e', \varepsilon_2 f'\}$ for some $\varepsilon_1, \varepsilon_2 \in \{-1, 1\}$. This implies $\gamma = \gamma'$ in view of (44), and concludes the proof. \square

The next lemma establishes weak uniqueness and stability properties for M -reduced bases, in the case of a strict inequality $2\mu(M) < \lambda_1(M)^2$.

Lemma 12. *Consider $M, M' \in S_2^+$, an M -reduced basis (e, f) , and an M' -reduced basis (e', f') . Let $\tau \geq 1$ be such that $\tau^{-2}M \leq M' \leq \tau^2 M$, in the sense of symmetric matrices. Assume either:*

(i) $2\mu(M) < \lambda_1(M)^2$, and $\tau = 1$ (i.e. $M' = M$).

(ii) $4\mu(M) \leq \lambda_1(M)^2$, and $\tau^4 \leq 1 + \frac{1}{3}\kappa(M)^{-2}$.

Then $\{e', f'\} \subset \{e, f, -e, -f\}$.

Proof. Denoting $\alpha := 2\mu(M)/\lambda_1(M)^2$, we obtain:

$$\begin{aligned} 4\langle e, M'f \rangle &= \|e + f\|_{M'}^2 - \|e - f\|_{M'}^2 \\ &\leq \tau^2 \|e + f\|_M^2 - \tau^{-2} \|e - f\|_M^2 \\ &= (\tau^2 - \tau^{-2})(\|e\|_M^2 + \|f\|_M^2) + 2(\tau^2 + \tau^{-2})\langle e, Mf \rangle \\ &\leq ((\tau^2 - \tau^{-2})(1 + \kappa(M)^2) + \alpha(\tau^2 + \tau^{-2}))\|e\|_M^2 \\ &\leq ((\tau^4 - 1)(1 + \kappa(M)^2) + \alpha(\tau^4 + 1))\|e\|_{M'}^2. \end{aligned}$$

In the fourth line we used Lemma 2, which implies that $\|f\|_M = \lambda_2(M) \leq \kappa(M)\lambda_1(M) = \kappa(M)\|e\|_M$, and Lemma 3 to bound $2\langle e, Mf \rangle$. Replacing α and τ with their assumed upper bounds, we obtain $2\langle e, M'f \rangle < \|e\|_{M'}^2$. Proceeding likewise, we obtain $2|\langle e, M'f \rangle| < \min\{\|e\|_{M'}^2, \|f\|_{M'}^2\}$. We may therefore apply Lemma 10 to M' and $(e, f, -e, -f)$, which implies $\{e', f'\} \subset \{e, f, -e, -f\}$ as announced. \square

3.3 Comparison of the stencils

We assume in this subsection that the scale parameter h is sufficiently small. Our assumption is stronger than the one used in §3.1, see (30), hence in particular there exists an Anisotropic Delaunay Triangulation \mathcal{T}_h . More precisely we assume that

$$\tau_h \leq \sqrt[4]{1 + 1/(3\kappa^2)} \quad \text{and} \quad \theta_h \leq \theta_0 := 1/(4\kappa). \quad (45)$$

See (26), (29), and Lemma 7 for the definition of κ, τ_h and θ_h respectively. For Lipschitz metrics, $\tau_h = 1 + \mathcal{O}(h)$ and $\theta_h = \mathcal{O}(h)$.

Our objective is to compare the stencils $V(p), V_h(p)$, of the AD-LBR (10) and of the ADT finite element discretization (17) respectively, at a point $p \in h\mathbb{Z}^2$. The next lemma shows that they are equal *unless* the lattice \mathbb{Z}^2 is almost orthogonal with respect to the local metric; a property quantified via $\mu(\mathbf{M}(p))$, see (40).

Lemma 13. *Let $p \in h\mathbb{Z}^2$, and let $M := \mathbf{M}(p)$. If $\mu(M) > \theta_h$, then $V_h(p) = V(p)$. In any case, one has for any M -reduced basis (e, f) :*

$$V_h(p) \supset \{e, f, -e, -f\} \quad (46)$$

$$V_h(p) \subset \{e, f, e + f, e - f, -e, -f, -e - f, f - e\} \quad (47)$$

Proof. We assume that $\langle e, Mf \rangle \leq 0$, up to replacing f with $-f$. Let $T \in \mathcal{T}_h$ be a triangle containing p , and let he_1, he_2, he_3 be the edges of T , oriented so that $e_1 + e_2 + e_3 = 0$. Using point (iii) of Lemma 7, and (27), we obtain for all $1 \leq i \leq 3$, with the convention $e_4 := e_1$

$$\langle e_i, Me_{i+1} \rangle \leq \theta_h \leq \theta_0 < \frac{1}{2\kappa} \leq \frac{1}{2} \min\{\|e_i\|_M^2, \|e_{i+1}\|_M^2\}. \quad (48)$$

Denote $E := \{e_1, e_2, e_3\}$, and $-E := \{-e_1, -e_2, -e_3\}$. Applying Lemma 10 to M and the points $(e_1, -e_3, e_2, -e_1, e_3, -e_2)$, we obtain that $\{e, f\} \subset E \cup (-E)$. Up to exchanging E with $-E$, we thus have $\{e, f\} \subset E$ or $\{e, -f\} \subset E$. Since the elements of E sum to zero, we conclude that

$$E = \{e, f, -e - f\} \quad \text{or} \quad E = \{e, -f, -e + f\}, \quad (49)$$

which implies (47).

If $\mu(M) = |\langle e, Mf \rangle| > \theta_h$, then (48) forbids the second case in (49). Thus $E = \{e, f, -e - f\}$, and therefore $V_h(p) \subset V(p)$, using (39).

Let $T \in \mathcal{T}_h$ be a triangle containing p and intersecting the half line $L := \{p + re; r > 0\}$. We know (49) that he is a vector edge of T (i.e. the difference between two vertices of T). The corresponding edge segment must be $[p, p + he]$, since otherwise $T \cap L$ would be empty. Thus $e \in V_h(p)$. Applying the same argument to $-e, f, -f$, we obtain (46).

If $\mu(M) > \theta_h$, then $h(e + f)$ is also a vector edge of any triangle $T \in \mathcal{T}_h$ containing p , since we eliminated

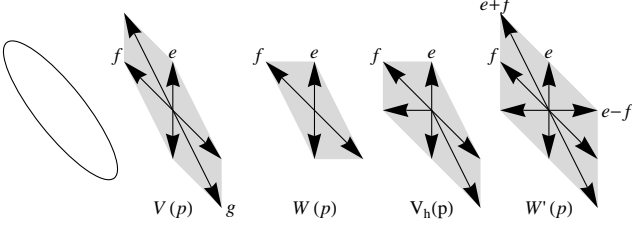


Figure 6: Consider a point $p \in h\mathbb{Z}^2$, and denote $M := \mathbf{M}(p)$. From left to right: ellipse $\{\|z\|_M \leq 1\}$, AD-LBR stencil $V(p)$, stencils $W(p) \subset V_h(p) \subset W'(p)$. For $W(p)$ and $W'(p)$ we assumed that $\mu(M) < \theta_0$, otherwise they are equal to $V(p)$. Note that $V(p)$, $W(p)$, $W'(p)$ only depend on M , while $V_h(p)$ depends on the structure of the triangulation \mathcal{T}_h .

the second case in (49). Reasoning as above we find that $\{e+f, -e-f\} \subset V_h(p)$, and therefore $V(p) \subset V_h(p)$. Thus $V(p) = V_h(p)$. This concludes the proof. \square

We introduce new stencils $W(p), W'(p)$, for $p \in \mathbb{R}^2$, defined as follows. Let $M := \mathbf{M}(p)$. If $\mu(M) \leq \theta_0$, then denoting by (e, f) an M -reduced basis,

$$W(p) := \{e, f, -e, -f\}, \quad (50)$$

$$W'(p) := \{e, f, e+f, e-f, -e, -f, -e-f, f-e\}. \quad (51)$$

On the other hand, if $\mu(M) > \theta_0$, then

$$W(p) := V(p) =: W'(p). \quad (52)$$

The previous lemma implies that $W(p) \subset V_h(p) \subset W'(p)$ for any $p \in h\mathbb{Z}^2$.

Lemma 14. *The stencils $W(p), W'(p)$, do not depend on the choice of $\mathbf{M}(p)$ -reduced basis.*

Proof. Let $M := \mathbf{M}(p)$. If $\mu(M) > \theta_0$, then $W(p), W'(p)$ are defined by (52), hence there is nothing to prove. Otherwise we obtain $\mu(M) \leq \theta_0 \leq 1/(4\kappa) \leq \lambda_1(M)^2/4$. Hence, by Lemma 12, any two M -reduced bases (e, f) , (e', f') , need to satisfy $\{e', f'\} \subset \{e, f, -e, -f\}$. In view of (50) and (51), they thus yield the same stencils $W(p), W'(p)$. \square

Let $\mathcal{F}_h, \mathcal{F}'_h$ be the energies associated to the stencils W, W' : for $u \in L^2(\Omega_h)$, extended to $h\mathbb{Z}^2$ by periodicity,

$$\begin{aligned} \mathcal{F}_h(u) &:= \sum_{z \in \Omega_h} \sum_{g \in W(z)} |u(z+hg) - u(z)|^2, \\ \mathcal{F}'_h(u) &:= \sum_{z \in \Omega_h} \sum_{g \in W'(z)} |u(z+hg) - u(z)|^2. \end{aligned}$$

The outline of the proof of Theorem 1 is as follows. We prove in Lemmas 16, 17 and 15 respectively that for any

$u \in L^2(\Omega_h)$:

$$|\mathcal{E}'_h(u) - \mathcal{E}_h(u)| \leq (\varepsilon_h + C_0\theta_h)\mathcal{F}'_h(u) \quad (53)$$

$$\mathcal{F}'_h(u) \leq C_1\mathcal{F}_h(u) \quad (54)$$

$$\mathcal{F}_h(u) \leq C_2\mathcal{E}_h(u), \quad (55)$$

where the constants C_0, C_1, C_2 only depend on the metric \mathbf{M} . Combining these inequalities, and recalling that $\theta_h = \mathcal{O}(h)$ and $\varepsilon_h = \mathcal{O}(h)$ for Lipschitz metrics (ε_h is defined in Lemma 9), we obtain

$$|\mathcal{E}'_h(u) - \mathcal{E}_h(u)| \leq ch\mathcal{E}_h(u),$$

for some constant $c = c(\mathbf{M})$. This establishes (13), and concludes the proof of Theorem 1.

For each $p \in \mathbb{R}^2$, we denote by $\eta_p, \eta'_p : \mathbb{Z}^2 \rightarrow \{0, 1\}$, the characteristic functions of $W(p)$ and $W'(p)$ respectively. The proofs of (53) and (55) immediately result from the comparison, in Lemmas 16 and 15 respectively, of the coefficients $\gamma_p, \gamma_p^h, \eta_p, \eta'_p$ appearing in the expressions of $\mathcal{E}_h, \mathcal{E}'_h, \mathcal{F}_h, \mathcal{F}'_h$.

In the following, it will be convenient to express the AD-LBR weights, and others, in terms of the scalar product associated to the Riemannian metric. We thus recall (8): for any $z \in \Omega$, and any $e, f \in \mathbb{R}^2$,

$$\langle e^\perp, \mathbf{D}(z)f^\perp \rangle = \mathbf{d}(z)\langle e, \mathbf{M}(z)f \rangle.$$

We also define the bounds $(0 < \underline{\mathbf{d}} \leq \overline{\mathbf{d}} < \infty)$

$$\underline{\mathbf{d}} := \min_{z \in \Omega} \mathbf{d}(z), \quad \overline{\mathbf{d}} := \max_{z \in \Omega} \mathbf{d}(z).$$

Lemma 15. *For any $p \in \mathbb{R}^2$, one has on \mathbb{Z}^2*

$$\eta_p \leq C_2\gamma_p, \quad \text{with } C_2 := 2\overline{\mathbf{d}}/\theta_0.$$

Proof. Let $M := \mathbf{M}(p)$, and let (e, f, g) be an M -obtuse superbase of \mathbb{Z}^2 . We can assume, thanks to Corollary 2, that (e, f) is an M -reduced basis. Then using (22)

$$2\mathbf{d}(p)\gamma_p(\pm f) \geq \frac{1}{2}\|e\|_M^2 \geq \frac{1}{2\kappa} = 2\theta_0,$$

hence $\gamma_p(\pm f) \geq \theta_0/\overline{\mathbf{d}}$, and likewise $\gamma_p(\pm e) \geq \theta_0/\overline{\mathbf{d}}$. If $\mu(M) \leq \theta_0$, then $W(p) = \{e, f, -e, -f\}$, and this concludes the proof.

Assume now that $\mu(M) > \theta_0$. Then

$$2\mathbf{d}(p)\gamma_p(\pm g) = -\langle e, Mf \rangle = \mu(M) \geq \theta_0,$$

hence $\gamma_p(\pm g) \geq \theta_0/(2\overline{\mathbf{d}})$. The result follows since $W(p) = \{e, f, g, -e, -f, -g\}$. \square

Let $p \in h\mathbb{Z}^2$ and let e_1, \dots, e_k be the consecutive elements of $V_h(p)$, in trigonometric order. We define for all $1 \leq i \leq k$, denoting $M := \mathbf{M}(p)$,

$$\tilde{\gamma}_p^h(e_i) := -\frac{\mathbf{d}(p)}{4}(\langle e_i - e_{i-1}, Me_{i-1} \rangle + \langle e_i - e_{i+1}, Me_{i+1} \rangle),$$

with the periodic conventions $e_{k+1} := e_1$, $e_0 := e_k$. We also set $\tilde{\gamma}_p^h = 0$ on $\mathbb{Z}^2 \setminus \{e_1, \dots, e_k\}$.

Lemma 16. *For any $p \in h\mathbb{Z}^2$, one has on \mathbb{Z}^2*

$$|\gamma_p^h - \tilde{\gamma}_p^h| \leq \varepsilon_h \eta'_p, \quad \text{and} \quad |\tilde{\gamma}_p^h - \gamma_p| \leq C_0 \theta_h \eta'_p, \quad (56)$$

where ε_h is given in Lemma 9, and $C_0 = 1/\underline{\mathbf{d}}$.

Proof. The coefficients γ_p , γ_p^h , $\tilde{\gamma}_p^h$, are all equal to zero outside of $W'(p)$. This holds by construction of γ_p , and by Lemma 13 for γ_p^h , $\tilde{\gamma}_p^h$. We may therefore forget about the presence of η'_p in (56).

First inequality. Lemma 9 states that $|\gamma_p^h - \tilde{\gamma}_p^h| \leq \varepsilon_h$ on \mathbb{Z}^2 , which concludes the proof.

Second inequality. If $\mu(M) > \theta_h$, then $V_h(p) = V(p)$. Comparing the definition of $\tilde{\gamma}_p^h$ with that of γ_p (11) we observe that $\tilde{\gamma}_p^h = \gamma_p$ on \mathbb{Z}^2 , which concludes the proof in this case.

Assume now that $\mu(M) \leq \theta_h$. Let (e, f, g) be an M -obtuse superbase of \mathbb{Z}^2 . We can assume, thanks to Corollary 2 that (e, f) is an M -reduced basis. Looking at (11) and denoting $\delta := 2\mathbf{d}(p)$, we find that

$$|\delta \gamma_p(\pm e) - \|f\|_M^2| = |\langle e, Mf \rangle| = \mu(M) \leq \theta_h.$$

Likewise $|\delta \gamma_p(\pm f) - \|e\|_M^2| \leq \theta_h$. In addition

$$\delta \gamma_p(\pm(e+f)) = \mu(M) \leq \theta_h, \text{ and } \gamma_p(\pm(e-f)) = 0.$$

Combining the definition of $\tilde{\gamma}_p^h$ with the description of the stencil $V_h(p)$ in Lemma 13, we obtain that

$$2\delta \tilde{\gamma}_p^h(e) = \left\{ \begin{array}{c} \langle f - e, Mf \rangle \\ \text{or} \\ \langle f + e, Mf \rangle \end{array} \right\} + \left\{ \begin{array}{c} \langle f - e, Mf \rangle \\ \text{or} \\ \langle f + e, Mf \rangle \end{array} \right\}.$$

In any case $|\delta \tilde{\gamma}_p^h(e) - \|f\|_M^2| \leq \theta_h$. The expressions and estimates of $\tilde{\gamma}_p^h$ at the points $-e, f, -f$ are obtained similarly. Likewise, using Lemma 13,

$$2\delta \tilde{\gamma}_p^h(e+f) = \left\{ \begin{array}{ll} \langle e, Mf \rangle + \langle e, Mf \rangle & \text{if } e+f \in V_h(p), \\ 0 & \text{otherwise.} \end{array} \right.$$

In any case $|\delta \tilde{\gamma}_p^h(e+f)| \leq \theta_h$. The expressions and estimates of $\tilde{\gamma}_p^h$ at the points $-(e+f), e-f, -(e-f)$ are similar. Comparing the above estimates of γ_p , $\tilde{\gamma}_p^h$, we obtain that $\delta|\gamma_p - \tilde{\gamma}_p^h| \leq 2\theta_h$ on $\{e, f, e+f, e-f, -e, -f, -e-f, f-e\} = W'(p)$. Since $\delta = 2\mathbf{d}(p) \geq 2\underline{\mathbf{d}} = 2/C_0$, this concludes the proof. \square

In the last lemma of this section, we control the contribution to the energy \mathcal{F}'_h of a stencil $W'(p)$, $p \in h\mathbb{Z}^2$, in terms of the contributions to \mathcal{F}_h of $W(p)$ and of the neighboring stencils $W(p+he)$, $e \in W(p)$. This leads to an estimate of \mathcal{F}'_h in terms of \mathcal{F}_h , which concludes the proof of Theorem 1.

Lemma 17. *One has $\mathcal{F}'_h(u) \leq C_1 \mathcal{F}_h(u)$, for any $u \in L^2(\Omega_h)$, with $C_1 := 17$.*

Proof. Consider a grid point $p \in h\mathbb{Z}^2$, and denote $M := \mathbf{M}(p)$. Assume first that $\mu(M) \leq \theta_0$, so that $W(p) \subsetneq W'(p)$. Consider also an arbitrary $g \in W'(p) \setminus W(p)$, and observe that $g = e + f$ for some M -reduced basis (e, f) .

We set $p' := p + e$ and $M' := \mathbf{M}(p')$. Applying point (ii) of Lemma 12, we find that (e, f) is also an M' -reduced basis. Indeed we have as required

$$4\mu(M) \leq 4\theta_0 = \kappa^{-1} \leq \lambda_1(M)^2,$$

using (27), and the assumption on τ follows from (45) and (28). Therefore

$$f \in W(p'), \text{ and } h^{-1}(p' - p) = e \in W(p'). \quad (57)$$

We obtain

$$\begin{aligned} |u(p+g) - u(p)|^2 &= |u(p+e+f) - u(p)|^2 \\ &\leq 2(|u(p+e+f) - u(p+e)|^2 + |u(p+e) - u(p)|^2) \\ &= 2(|u(p'+f) - u(p')|^2 + |u(p+e) - u(p)|^2). \end{aligned} \quad (58)$$

Denote, for all $q \in h\mathbb{Z}^2$,

$$\begin{aligned} \mathcal{F}_h(u; q) &:= \sum_{g \in W(q)} |u(q+hg) - u(q)|^2, \\ \mathcal{F}'_h(u; q) &:= \sum_{g \in W'(q)} |u(q+hg) - u(q)|^2. \end{aligned}$$

Using (58), we obtain

$$\mathcal{F}'_h(u; p) - \mathcal{F}_h(u; p) \leq \mathcal{G}_h(u; p) \quad (59)$$

where $\mathcal{G}_h(u; p)$ is given by

$$\left\{ \begin{array}{ll} 4\mathcal{F}_h(u; p) + 2 \sum_{g \in W(p)} \mathcal{F}_h(u; p+g), & \text{if } \mu(\mathbf{M}(p)) \leq \theta_0 \\ 0, & \text{if } \mu(\mathbf{M}(p)) > \theta_0 \end{array} \right.$$

When $\mathcal{F}_h(u; p')$ appears in $\mathcal{G}_h(u; p)$, with $p, p' \in h\mathbb{Z}^2$, $p \neq p'$, we have $h^{-1}(p' - p) \in W(p')$, see (57). For each $p' \in h\mathbb{Z}^2$, there are thus at most $\#(W(p')) \leq 6$ points $p \in h\mathbb{Z}^2 \setminus \{p'\}$ such that $\mathcal{F}_h(u; p')$ appears in $\mathcal{G}_h(u; p)$. Summing (59) over $p \in \Omega_h$, we thus obtain $\mathcal{F}'_h(u) - \mathcal{F}_h(u) \leq (4 + 2 \times 6)\mathcal{F}_h(u)$ (the constant could easily be improved), which concludes the proof. \square

4 Numerical experiments

We compare our scheme AD-LBR with a family of other schemes: finite difference, finite elements, and two schemes from the image processing literature. We begin with

a quantitative comparison for the discretization of the restoration equation, in a synthetic case where the exact solution is analytically available for reference. The second test case is a qualitative comparison of Coherence-Enhancing Diffusion (CED) [25], on a real image and the quality assessment is by visual inspection. Finally we present a 3D implementation of AD-LBR for proof of feasibility, featuring a synthetic CED experiment, and a application of Edge Enhancing Diffusion to MRI data.

4.1 The different schemes

Our two dimensional numerical experiments feature the following six numerical schemes for anisotropic diffusion.

- AD-LBR: the scheme presented in this work.
- Finite Differences (FD). The gradient and the divergence are discretized using standard centered finite differences [15], see Remark 3 for details. This approach, arguably the most straightforward, leads to a 9 point stencil.
- Bilinear Finite Elements (Q1). Bilinear finite elements, also referred to as Q1 finite elements, are linear with respect to each space direction. This amounts to use a 9 points stencil, where the coefficients are different from the previous scheme.
- Weickert-Scharr scheme (WS). This scheme, introduced in [27], is based on a second order approximation of the gradient using a 3×3 centered stencil. As a result, it offers good accuracy and rotation invariance when applied to sufficiently smooth functions, but lacks robustness guarantees such as the discrete maximum principle and spectral correctness (see §1), even for $\mathbf{D} = \text{Id}$. The stencil for this scheme has size 5×5 .
- Weickert's Non-Negative scheme (W-NN). The coefficients of this scheme, detailed in [25] page 95, are non-negative as long as the anisotropy ratio (9) satisfies $\kappa \leq 1 + \sqrt{2} \sim 2.41$.
- Axes-directed Non-Negative scheme (A-NN). This six point non-negative scheme is implicitly defined in the proof Theorem 6 in [25], and can be regarded as a generalisation of W-NN. See Remark 1 below for details. Among the 6 points of the stencil, 4 points are along the axes of coordinates.

Note that other schemes exist, see for instance [16, 28]. While an exhaustive comparison is in principle desirable, it could not be done here due to time and space constraints.

To fix the ideas and illustrate the difference between the schemes, we propose to compute the stencil and the coefficients for different *constant* diffusion tensors \mathbf{D} , in

isotropic and anisotropic cases. Denoting by R the matrix of rotation by the angle $\theta = \pi/6$, and by $\kappa \geq 1$ the chosen anisotropy ratio, we set, identically on \mathbb{R}^2 :

$$\mathbf{D} := R \begin{pmatrix} 1 & 0 \\ 0 & \kappa^{-2} \end{pmatrix} R^T. \quad (60)$$

The results are presented in Tables 1 and 2. Note that for the two last cases (anisotropy $\kappa = \sqrt{10}$ and $\kappa = \sqrt{50}$) the AD-LBR stencil contains points that are outside the 3×3 neighborhood of the pixel. However the stencil contains 6 points, as expected. This contrasts with the schemes FD, Q1, W-NN where only the 3×3 neighborhood is involved. Another observation is that the off-center stencil coefficients of the AD-LBR are non-positive (this gives non-negative off-diagonal coefficients for $\text{div}(\mathbf{D}\nabla)$), in contrast with schemes FD, Q1, WS, and with scheme W-NN for anisotropy $\kappa > 1 + \sqrt{2}$. This is an essential property of AD-LBR (and A-NN), and as a consequence our scheme satisfies, unconditionally, the discrete maximum principle [1, 6].

The largest eigenvalue of the discrete operator $-\text{div}(\mathbf{D}\nabla)$ is given in Table 3, for the different schemes. It turns out that AD-LBR has in most cases the smallest eigenvalues among all schemes, except for scheme WS and occasionally A-NN. This property allows (although this was not done in our numerical experiments) to use larger time steps for AD-LBR than for the other schemes, when solving parabolic equations (2) or (64) with an explicit time discretization.

Operator splitting is a classical approach to further increase the timestep in (potentially anisotropic and non-linear) diffusion PDEs [25, 26, 2]. The AD-LBR is compatible with Additive Operator Splitting, by applying Remark (e) page 111 in [25], although the efficiency of this technique is here compromised by the potentially large number of directions in our adaptive stencils. Let us also mention Multiplicative Operator Splittings, and Additive-Multiplicative Operator Splittings, which allow to combine different time-steps [2, 8]. None of these methods was used in our experiments.

Remark 1 (Axes-directed non negative six point scheme). *The following six point scheme A-NN is, in our belief, the best possible implementation of the constructive proof in [25] of the existence of non-negative schemes for two dimensional anisotropic diffusion. Like AD-LBR, this scheme is defined by the data at each point $z \in \Omega$ of a stencil $V(z)$, and of non-negative weights γ_z .*

Let $\mathbf{D}(z) = \begin{pmatrix} a & b \\ b & c \end{pmatrix}$. In the diagonal case $b = 0$, the scheme A-NN relies on the classical four points stencil. Otherwise note that

$$\frac{a}{|b|} - \frac{|b|}{c} = \frac{ac - b^2}{|b|c} > 0.$$

Let $p, q \in \mathbb{Z} \setminus \{0\}$ be such that $bpq \geq 0$,

$$\frac{|b|}{c} \leq \left| \frac{p}{q} \right| \leq \frac{a}{|b|}, \quad (61)$$

and $\max(|p|, |q|)$ is minimal. The scheme A-NN is defined by the six point stencil

$$V(z) := \{(\pm 1, 0), (0, \pm 1), \pm(p, q)\}$$

and the non-negative weights

$$\begin{aligned} 2\gamma(\pm 1, 0) &= a - \frac{p}{q}b, & 2\gamma(0, \pm 1) &= c - \frac{q}{p}b, \\ 2\gamma(\pm(p, q)) &= \frac{b}{pq}. \end{aligned}$$

These coefficients are non-negative by construction, and consistency (5) is easily checked. Contrary to AD-LBR, the coordinate axes play a privileged role in A-NN. This introduces axis aligned artifacts which are visible in Figure 12 (g).

Remark 2 (Stencil radius). The two dimensional stencils of AD-LBR coincide with those of FM-LBR, a numerical scheme for anisotropic static Hamilton-Jacobi PDEs introduced in [14] the second author. As shown in Proposition 1.6 of [14], the euclidean radius

$$r = \max\{\|v\|; v \in V(z)\}$$

of this stencil is bounded by $\kappa(\mathbf{D}(z))$.

In contrast, consider for $0 < \varepsilon < 1/4$ the matrix

$$D := \begin{pmatrix} 1 & 1 - 2\varepsilon \\ 1 - 2\varepsilon & 1 - 3\varepsilon \end{pmatrix}.$$

It follows from (61) that $1 + \varepsilon \leq p/q + \mathcal{O}(\varepsilon^2) \leq 1 + 2\varepsilon$. From this point, one easily obtains that $q \gtrsim \varepsilon^{-1} \approx \kappa(D)^2$. The radius of the A-NN stencil, at a point $z \in \Omega$, may thus be of the order of $\kappa(\mathbf{D}(z))^2$. The radii of the AD-LBR and A-NN stencils, computed for diffusion tensors of anisotropy $\kappa = 10$ and of various orientations, are illustrated on Figure 7.

Remark 3 (Scheme FD). The operator $\operatorname{div}(\mathbf{D} \nabla \cdot)$ is discretized using centered finite differences [15]. This involves quantities defined at half integer indices, and in particular the diffusion tensor is here given on the offsetted grid $(i + 1/2, j + 1/2)$, $(i, j) \in \mathbb{Z}^2$. For the sake of readability, we thus define $i^+ := i + 1/2$ and $i^- := i - 1/2$. The gradient operator is discretized by:

$$(\partial_x u)_{i^+, j} = u_{i+1, j} - u_{i, j}, \quad (\partial_y u)_{i, j^+} = u_{i, j+1} - u_{i, j}.$$

The divergence is defined as follows:

$$\begin{aligned} \operatorname{div}(\mathbf{D} \nabla u)_{i, j} &= \partial_x(\mathbf{D}^{11} \partial_x u + \mathbf{D}^{12} \partial_y u)_{i, j} \\ &\quad + \partial_y(\mathbf{D}^{21} \partial_x u + \mathbf{D}^{22} \partial_y u)_{i, j}, \end{aligned}$$

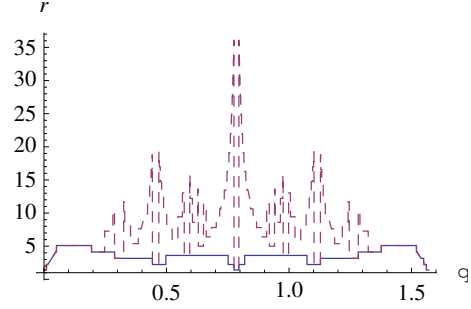


Figure 7: Radius of the AD-LBR stencil (plain), and of the A-NN stencil (dashed), for a matrix D_θ of anisotropy ratio $\kappa = 10$ and eigenvector $(\cos \theta, \sin \theta)$. The AD-LBR stencil is here always the smallest, and its radius does not exceed 5.1, versus 36.1 for A-NN.

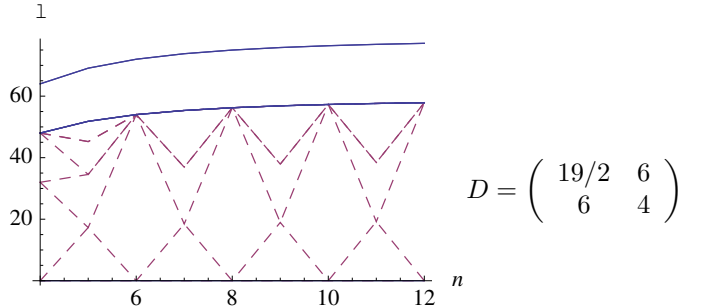


Figure 8: The seven smallest eigenvalues of the operator $-\operatorname{div}(D \nabla)$, on $[0, 1]^2$ with periodic boundary conditions, discretized on a $n \times n$ grid with $4 \leq n \leq 12$. Plain: AD-LBR discretization; Dashed: A-NN discretization. Some eigenvalues have multiplicities, hence less than 14 graphs are visible. Eigenvalues of the A-NN discretization here “oscillate” with the dimension. In contrast, thanks to its asymptotic equivalence with a finite element scheme, the smallest eigenvalues of the AD-LBR discretization converge towards those of the continuous partial differential operator.

with

$$\begin{aligned}
(\mathbf{D}^{11}\partial_x u)_{i^+,j} &= \frac{1}{2} \left(\mathbf{D}_{i^+,j^+}^{11} + \mathbf{D}_{i^+,j^-}^{11} \right) (\partial_x u)_{i^+,j}, \\
\partial_x (\mathbf{D}^{11}\partial_x u)_{i,j} &= (\mathbf{D}^{11}\partial_x u)_{i^+,j} - (\mathbf{D}^{11}\partial_x u)_{i^-,j}, \\
(\mathbf{D}^{21}\partial_x u)_{i^+,j^+} &= \frac{1}{2} \mathbf{D}_{i^+,j^+}^{21} ((\partial_x u)_{i^+,j} + (\partial_x u)_{i^+,j+1}), \\
\partial_y (\mathbf{D}^{21}\partial_x u)_{i,j} &= \frac{1}{2} ((\mathbf{D}^{21}\partial_x u)_{i^+,j^+} - (\mathbf{D}^{21}\partial_x u)_{i^+,j^-} \\
&\quad + (\mathbf{D}^{21}\partial_x u)_{i^-,j^+} - (\mathbf{D}^{21}\partial_x u)_{i^-,j^-}),
\end{aligned}$$

and similar terms involving $\partial_y u$.

Table 1: The stencil coefficients for different *constant* diffusion tensors, and the different schemes presented. The value of the anisotropy ratio κ is given in the second row, and the orientation of the principal axis is $\theta = \pi/6$, see (60). The bold coefficient indicates the center node. In some examples we present for clarity reasons only half of the stencil (the other half can be deduced by symmetry). Stencil entries are highlighted when they are positive and off-center - an undesirable property which gives rise to stability issues. For small anisotropies, $\kappa \leq 1 + \sqrt{2}$, one has AD-LBR = W-NN = A-NN.

| κ | $\kappa = 1$ ($\mathbf{D} = \text{Id}$) | $\kappa = \sqrt{2}$ |
|--------------------|--|---|
| stencil for AD-LBR | $\begin{matrix} 0 & -1 & 0 \\ -1 & \mathbf{4} & -1 \\ 0 & -1 & 0 \end{matrix}$ | $\begin{matrix} 0 & -0.41 & -0.22 \\ -0.66 & \mathbf{2.57} & -0.66 \\ -0.22 & -0.41 & 0 \end{matrix}$ |
| stencil for FD | $\begin{matrix} 0 & -1 & 0 \\ -1 & \mathbf{4} & -1 \\ 0 & -1 & 0 \end{matrix}$ | $\begin{matrix} \mathbf{0.11} & -0.63 & -0.11 \\ -0.88 & \mathbf{3} & -0.88 \\ -0.11 & -0.63 & \mathbf{0.11} \end{matrix}$ |
| stencil for Q1 | $\frac{1}{3} \begin{pmatrix} -1 & -1 & -1 \\ -1 & \mathbf{8} & -1 \\ -1 & -1 & -1 \end{pmatrix}$ | $\begin{matrix} -0.14 & -0.13 & -0.36 \\ -0.38 & \mathbf{2} & -0.38 \\ -0.36 & -0.13 & -0.14 \end{matrix}$ |
| stencil for WS | $\begin{matrix} -0.1 & -0.06 & -0.02 \\ \mathbf{0.12} & 0 & -0.06 \\ \mathbf{0.46} & \mathbf{0.12} & -0.1 \\ \mathbf{0.12} & 0 & -0.06 \\ -0.1 & -0.06 & -0.02 \end{matrix}$ | $\begin{matrix} -0.06 & -0.05 & -0.02 \\ \mathbf{0.01} & -0.04 & -0.06 \\ \mathbf{0.35} & \mathbf{0.07} & -0.09 \\ \mathbf{0.01} & \mathbf{0.04} & -0.04 \\ -0.06 & -0.02 & -0.01 \end{matrix}$ |
| stencil for W-NN | $\begin{matrix} 0 & -1 & 0 \\ -1 & \mathbf{4} & -1 \\ 0 & -1 & 0 \end{matrix}$ | $\begin{matrix} 0 & -0.41 & -0.22 \\ -0.66 & \mathbf{2.57} & -0.66 \\ -0.22 & -0.41 & 0 \end{matrix}$ |
| stencil for A-NN | $\begin{matrix} 0 & -1 & 0 \\ -1 & \mathbf{4} & -1 \\ 0 & -1 & 0 \end{matrix}$ | $\begin{matrix} 0 & -0.41 & -0.22 \\ -0.66 & \mathbf{2.57} & -0.66 \\ -0.22 & -0.41 & 0 \end{matrix}$ |

4.2 A test case with an explicit solution

Consider an image $v \in L^2(\Omega)$, defined on a domain Ω , and a diffusion tensor field $\mathbf{D} : \Omega \rightarrow S_2^+$. A classical approach to restore the image v , if it has been corrupted by additive noise, is to find $u \in H^1(\Omega)$ which minimizes:

$$j(u) = \int_{\Omega} |u - v|^2 + \lambda \int_{\Omega} \|\nabla u\|_{\mathbf{D}}^2. \quad (62)$$

Table 2: The stencil coefficients for different metrics and the different schemes presented, similarly to Table 1 but with more pronounced anisotropies. For the scheme A-NN some points of the stencil are too far from the center node to be represented here, so we indicate the coordinates of these points and the associated coefficient.

| κ | $\kappa = \sqrt{10}$ | | | $\kappa = \sqrt{50}$ | | |
|--------------------|------------------------|--------------|-------------|------------------------|-------------|-------------|
| stencil for AD-LBR | 0 | -0.26 | -0.06 | 0 | -0.11 | -0.16 |
| | 1.16 | -0.26 | 0 | 0.55 | -0.01 | 0 |
| stencil for FD | 0.19 | -0.32 | -0.19 | 0.21 | -0.27 | -0.21 |
| | -0.77 | 2.2 | -0.77 | -0.76 | 2.04 | -0.76 |
| | -0.19 | -0.32 | 0.19 | -0.21 | -0.27 | 0.21 |
| stencil for Q1 | 0.01 | 0.04 | -0.38 | 0.04 | 0.08 | -0.38 |
| | -0.41 | 1.47 | -0.41 | -0.42 | 1.36 | -0.42 |
| | -0.38 | 0.04 | 0.01 | -0.38 | 0.08 | 0.04 |
| stencil for WS | -0.02 | -0.04 | -0.02 | -0.02 | -0.04 | -0.02 |
| | 0.09 | -0.08 | -0.07 | 0.09 | -0.08 | -0.07 |
| | 0.25 | 0.04 | -0.08 | 0.24 | 0.03 | -0.08 |
| | 0.09 | 0.08 | -0.02 | 0.09 | 0.08 | -0.02 |
| | -0.02 | 0.004 | -0.003 | -0.02 | 0.01 | -0.002 |
| stencil for W-NN | 0 | 0.06 | -0.39 | 0 | 0.16 | -0.42 |
| | -0.39 | 1.42 | -0.39 | -0.33 | 1.19 | -0.33 |
| | -0.39 | 0.06 | 0 | -0.42 | 0.16 | 0 |
| stencil for A-NN | -0.07 | 0 | | -0.01 | 0 | |
| | 0.64 | -0.19 | | 0.17 | -0.05 | |
| | $\gamma(3, 2) = -0.06$ | | | $\gamma(5, 3) = -0.03$ | | |

In other words, u is a penalized least squares approximation of v . The parameter $\lambda > 0$ should be adjusted so as to avoid excessive smoothing (for large λ), or insufficient denoising (for small λ). The solution u can be characterized as the solution to the static elliptic PDE:

$$\begin{cases} -\lambda \operatorname{div}(\mathbf{D} \nabla u) + u = v, & \text{on } \Omega. \\ \langle \nabla u, n \rangle = 0, & \text{on } \partial \Omega. \end{cases} \quad (63)$$

In applications [20, 24] the diffusion tensor \mathbf{D} is usually adapted to the local image structure, in order to avoid smoothing the edges of v . We construct below a test case (image v and tensor field \mathbf{D}), for which the solution u is known analytically.

In order to obtain an analytic solution, we first consider a separable problem where the image is invariant by translation along the horizontal axis, and the metric is constant with axes parallel to the coordinate axes. This first problem is invariant under translations along the x -axis, and therefore boils down to a 1-dimensional problem. This separable problem is then transported by a diffeomorphism in order to obtain a new problem where the axes of the metric are no more parallel to the coordinate axes.

The analytical image is composed of a black and a white stripe: $v_0(x, y) = \mathbf{1}_{y < 0.5}$, see Figure 9. Given $\kappa \geq 1$, we consider the constant diffusion tensor

$$\mathbf{D}_0 = \begin{pmatrix} 1 & 0 \\ 0 & \kappa^{-2} \end{pmatrix}.$$

Table 3: Largest eigenvalue of the discretized operator $-\text{div}(\mathbf{D}\nabla)$, for the constant metric $\mathbf{D} = D$, where the matrix D is given on Tables 1 and 2. The time step, in the explicit discretization of (64), should not exceed the inverse of this value.

| κ | $\kappa = 1$ | $\kappa = \sqrt{2}$ | $\kappa = \sqrt{10}$ | $\kappa = \sqrt{50}$ |
|-------------------|--------------|---------------------|----------------------|----------------------|
| eigenvalue AD-LBR | 8 | 4.27 | 2.06 | 1.06 |
| eigenvalue FD | 8 | 6.22 | 5.06 | 4.85 |
| eigenvalue Q1 | 5.7 | 4.94 | 4.32 | 4.20 |
| eigenvalue WS | 1 | 1 | 1 | 1 |
| eigenvalue W-NN | 8 | 4.27 | 3.1 | 3.02 |
| eigenvalue A-NN | 8 | 4.27 | 1.04 | 0.3 |

The analytical solution u_0 of (63), applied to \mathbf{D}_0 and v_0 , is known in the case of the infinite domain $\Omega = \mathbb{R}^2$. In Fourier domain all the coefficients are real and:

$$\widehat{u}_0(\xi) = \widehat{v}_0(\xi)/(1 + \langle \xi, \mathbf{D}_0 \xi \rangle).$$

This separable problem is transformed using the following diffeomorphism: for $(x, y) \in \Omega$

$$f(x, y) = (x, y + \alpha \cos(2\pi x)).$$

The Jacobian of f is

$$J(x, y) = \begin{pmatrix} 1 & 0 \\ -2\pi\alpha \sin(2\pi x) & 1 \end{pmatrix}$$

We apply the different restoration schemes to the image $v = v_0 \circ f$, and the following diffusion tensor:

$$\begin{aligned} \mathbf{D}(z) &= |\det J(z)| J(z)^{-1} \mathbf{D}_0 (J(z)^{-1})^T \\ &= J(z)^{-1} \mathbf{D}_0 (J(z)^{-1})^T = \begin{pmatrix} 1 & s \\ s & s^2 + \kappa^{-2} \end{pmatrix}, \end{aligned}$$

where we denoted $z = (x, y) \in \Omega$ and $s = 2\pi\alpha \sin(2\pi x)$. The numerical solution is compared to the analytical function $u = u_0 \circ f$, which is the exact solution in the case of the infinite domain $\Omega = \mathbb{R}^2$. This numerical solution was obtained on the bounded domain $\Omega = [0, 1]^2$, equipped with reflecting boundary conditions. Numerical evidence suggests that this change of domain and of boundary conditions has only an anecdotic impact on the solution of (63), with the parameters chosen in this test case.

We used $\alpha := 1/3$ in the numerical experiments. The maximum value of $\kappa(\mathbf{D}(x))$, among all $x \in \Omega$, is equivalent to $\kappa_{\max} := \kappa\sqrt{1 + (2\pi\alpha)^2} \simeq 2.3\kappa$.

4.3 Results for the synthetic test case

We present in Figure 10 the performance results of the different schemes, for different values of the anisotropy κ ,

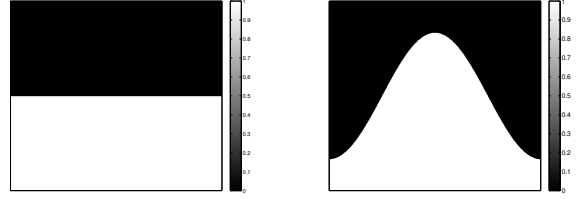


Figure 9: Left: image v_{ana} . Right: image $v = v_{\text{ana}} \circ f$ transformed by the diffeomorphism f .

obtained on a series of grids of size ranging from 100×100 to 1200×1200 . The anisotropy varies from $\kappa = 2$ to $\kappa = 10$, which are relevant values for imaging applications, see the numerical experiments in §4.4. The quality of a scheme is measured by the L^2 difference and the H^1 semi-norm difference between the numerical solution and the analytical solution. Note that the error is concentrated close to the discontinuity, since the solution tends rapidly to a constant (0 or 1) far from the discontinuity. We chose the smoothing parameter $\lambda = 10^{-3}$ in (62). The linear equation obtained by the discretization of (63) is solved using Conjugate Gradient.

We also tested extreme anisotropies, $\kappa \geq 100$ (thus $\kappa_{\max} \geq 230$), which can be relevant in physics related applications. None of the tested schemes showed convincing results: methods based on fixed stencils fail because the discrete operator loses positivity, while the AD-LBR (and A-NN even more) suffers from under-sampling due to the large radius of its stencils. We thus refer to [7] for a radically different approach tailored for this setting. This method introduces an auxiliary one-dimensional unknown, which is constant on the field lines (obtained in a preprocessing step) of the anisotropy direction field, and varies orthogonally to them.

The performance advantage of the AD-LBR is particularly clear when the error is measured in the H^1 semi-norm: for the anisotropy $\kappa = 10$ and the resolution 500×500 , which are relevant values in image processing, AD-LBR outperforms its alternatives by a factor ranging from 3 to 5.

4.4 Coherence-enhancing diffusion

In order to document the interest of our discretization, we implement Coherence-Enhancing Diffusion [25] using the different numerical schemes at our disposal. The following parabolic equation is considered:

$$\partial_t u = \text{div}(\mathbf{D}(J_\rho(\nabla u_\sigma)) \nabla u). \quad (64)$$

This equation is non-linear since the diffusion tensor depends on the solution u . This tensor also depends on four user defined parameters $\sigma, \rho, C \in \mathbb{R}_+$, $\alpha \in]0, 1[$. Let K_σ

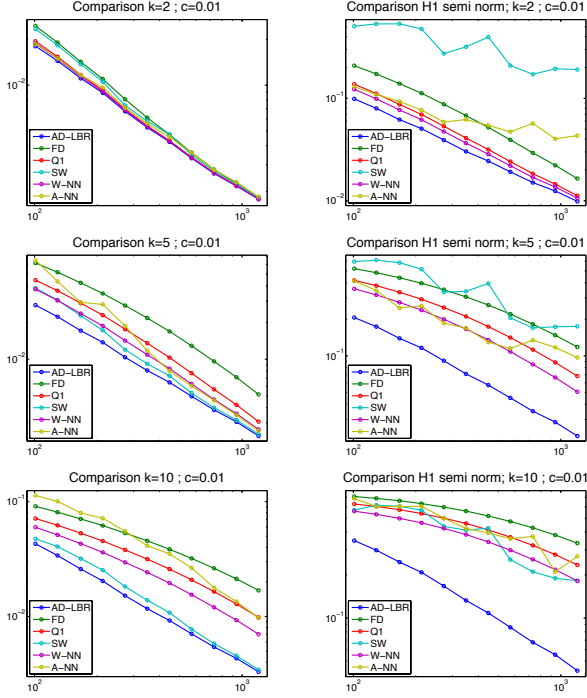


Figure 10: Numerical results for the synthetic test case, with different values of the anisotropy factor: $\kappa = 2$ (top row), $\kappa = 5$ (middle row), $\kappa = 10$ (bottom row). Vertical axis: relative error in L^2 norm (left column), or H^1 semi-norm (right column), for the six schemes tested. Horizontal axis: integer N , where the image resolution is $N \times N$. Since the tested schemes are first order, numerical error is expected to be proportional to N^{-1} . Log-log scale.

(resp. K_ρ), be the Gaussian kernel of variance σ (resp. ρ). Define the convolution $u_\sigma := K_\sigma \star u$, and the structure tensor $J_\rho := K_\rho \star (\nabla u_\sigma \nabla u_\sigma^T)$. The diffusion tensor $\mathbf{D}(J_\rho)$ possesses the same eigenvectors (v_1, v_2) as J_ρ , and if the eigenvalues of J_ρ are $\mu_1 \geq \mu_2$ then the eigenvalues of $\mathbf{D}(J_\rho)$ are

$$\begin{aligned} \lambda_1 &:= \alpha \\ \lambda_2 &:= \alpha + (1 - \alpha) \exp\left(\frac{-C}{(\mu_1 - \mu_2)^2}\right). \end{aligned}$$

This ensures that one smoothes preferably along the coherence direction v_2 , with a diffusivity that increases with respect to the coherence $(\mu_1 - \mu_2)^2$. When the time parameter t becomes large, the image tends to a constant image, therefore it is necessary to stop the process at some finite time T . The ratio of the eigenvalues is bounded by $\lambda_2/\lambda_1 \leq 1/\alpha$, hence $\kappa \leq 1/\sqrt{\alpha}$.

We used an explicit time discretization for (64), with time step Δt . The image u^{n+1} at time $(n+1)\Delta t$ is defined

by the explicit equation:

$$\frac{u^{n+1} - u^n}{\Delta t} = \text{div}(\mathbf{D}(J_\rho(\nabla u_\sigma^n)) \nabla u^n).$$

The parameters used in our simulation were: $\sigma = 0.5$, $\rho = 4$, $C = 10^{-5}$, $\alpha = 10^{-2}$ and $\Delta t = 0.02$. This gives a maximum anisotropy of $\kappa = 10$. The algorithm was applied to a fingerprint image. The results obtained for $T = 10$ are shown in Figures 11 and 12, and they document the ability of our scheme to close interrupted lines more efficiently than the other schemes. The largest eigenvalue of the discrete operator $-\text{div}(\mathbf{D}\nabla)$ at $t = 0$ is given in Table 4 for the different schemes. As was already noticed in the constant metric case, it turns out that AD-LBR has the smallest eigenvalues among all schemes, except for scheme WS. This property allows (although this was not done in our numerical experiments) to use larger time steps for AD-LBR than for the other schemes.

Note also that ridges are clearer, and valleys are darker, using AD-LBR than with the other schemes. (Gray-scale range is the same for all images, see also Figure 13). This reflects the fact that AD-LBR avoids, better than the other schemes, smoothing transversally to the orientation encoded in the continuous anisotropic PDE (64).

Remark 4 (Computation time). *Numerical solvers of the parabolic PDE (64) combine three main components: (i) Constructing the diffusion tensor. (ii) Assembling the discretization stencils and the operator sparse matrix. (iii) Performing an explicit time step. Components (i) and (ii) are executed exactly the same number of times, while step (iii) is generally more frequent: in order to save CPU time, one typically does not update the diffusion operator at each time step. We produced a C++ implementation of AD-LBR, within the Insight Toolkit open source library. Although our code is neither parallel nor aggressively optimized, we believe that comparing the CPU times for steps (i), (ii) and (iii) is informative, and allows to estimate the additional cost of AD-LBR which is essentially contained in step (ii).*

For our 2D Coherence-Enhancing Diffusion (CED) experiment, on the 512×512 fingerprint image, (i) takes 0.21s, (ii) 0.027s, (iii) 0.005s. For our 3D CED Experiment, on $100 \times 100 \times 100$ synthetic data, (i) takes 1.35s, (ii) 0.51s, (iii) 0.035s. In both cases, the AD-LBR specific step (ii) is dominated by the construction of the diffusion tensor (i). Step (ii) may also be dominated by the mere cost (iii) of iterations, provided the operator is updated less than once every 6 explicit steps in 2D (14 in 3D). To our eyes, the limited additional cost (ii) of AD-LBR is acceptable in view of the strong theoretical guarantees, and qualitative improvements, brought by this scheme.

Table 4: Largest eigenvalue of the discretized operator $-\operatorname{div}(\mathbf{D}\nabla)$, where $\mathbf{D} = \mathbf{D}(J_\rho(\nabla u_\sigma))$ at $t = 0$.

| scheme | AD-LBR | FD | Q1 | WS | W-NN | A-NN |
|------------|--------|------|------|------|------|------|
| eigenvalue | 3.75 | 5.67 | 5.09 | 0.96 | 3.83 | 6.23 |

4.5 3-dimensional experiments

In order to illustrate the feasibility of our scheme in 3D space, we present the action of anisotropic diffusion PDEs on two examples. The first example is a 3D analog of the synthetic test case presented in [27], featuring Coherence-Enhancing Diffusion. The second one is the application of Edge-Enhancing Diffusion to a MRI scan.

Synthetic example

The original, radially varying image is defined on the cube $[0, 1]^3$. The gray-level at a point x is defined by

$$u^0(x) = \cos(2(r/R)^3),$$

where $r := |x|$ and $R := 1/2$. This image presents a series of concentric level-sets. We present in Figure 14 the level sets $\{u^0 = 0\}$, and a slice through the plane $z = 0.7$.

The image u^0 is perturbed by

$$u := u^0 + n,$$

where n is an additive Gaussian noise of variance $\sigma = 0.5$. The reconstructed image is obtained using a 3D Coherence-Enhancing Diffusion PDE [25], similar to the 2D one in section 4.4:

$$\partial_t u = \operatorname{div}(\mathbf{D}(J_\rho(\nabla u_\sigma))\nabla u),$$

where J_ρ is the structure tensor defined by $J_\rho := K_\rho \star (\nabla u_\sigma \nabla u_\sigma^T)$, $u_\sigma := K_\sigma \star u$. The tensor $\mathbf{D}(J_\rho)$ possesses the same eigenvectors (v_1, v_2, v_3) as J_ρ , and if the eigenvalues of J_ρ are $\mu_1 \geq \mu_2 \geq \mu_3$ then the eigenvalues of $\mathbf{D}(J_\rho)$ are

$$\begin{aligned} \lambda_1 &:= \alpha \\ \lambda_2 &:= \alpha + (1 - \alpha) \exp\left(\frac{-C}{(\mu_1 - \mu_2)^2}\right), \\ \lambda_3 &:= \alpha + (1 - \alpha) \exp\left(\frac{-C}{(\mu_1 - \mu_3)^2}\right), \end{aligned}$$

where $\alpha = 10^{-2}$. The anisotropy ratio is bounded by $\kappa = 1/\sqrt{\alpha} = 10$. We used the values $\sigma = 0.5$, $\rho = 4$. The problem is discretized using 100^3 voxels. We present in Figure 14 the noisy image u (levelset 0 and planar slice) and the result after 20 time-steps of $\Delta t = 10^{-3}$.

3D MRI data

The data is a $256 \times 256 \times 100$ Magnetic Resonance Imaging scan of a skull, and was obtained from the "University of North Carolina Volume Rendering Test Data Set" archive.

The reconstructed image is obtained using a 3D Edge-Enhancing Diffusion PDE [25], which differs from the above Coherence-Enhancing Diffusion one by the choice of the diffusion tensor eigenvalues. The optimal choice of these eigenvalues indeed depends on the application, and is still an active subject of research [13]. With the above notations, the eigenvalues of $\mathbf{D}(J_\rho)$ are

$$\begin{aligned} \lambda_1 &:= 1 - \exp\left(\frac{-C}{\mu_1^2}\right) \\ \lambda_2 &:= 1 - \exp\left(\frac{-C}{\mu_2^2}\right), \\ \lambda_3 &:= 1. \end{aligned}$$

We used the values $\sigma = 0.5$, $\rho = 4$. In our experiment, the maximum anisotropy ratio was $\kappa = 11.2$. We present in Figure 15 the original image and two slices of the result after 10 time-steps of $\Delta t = 10^{-4}$.

Conclusion

We introduced in this paper a new numerical scheme, AD-LBR, for anisotropic diffusion in image processing. This scheme is non-negative, and its stencils have a limited support: 6 points in 2D, 12 points in 3D. The former property implies that our scheme respects the maximum principle of Alvarez, Guichard, Lions and Morel, which is an essential feature of parabolic PDEs.

AD-LBR outperformed all tested alternatives in a quantitative numerical experiment: a test case in which approximate numerical solutions are compared against a known analytical solution. In a second qualitative test case, different schemes were used to enhance a fingerprint image. Our scheme appears here to close more efficiently the lines of the fingerprint, and to diffuse less orthogonally to the lines. This is precisely the purpose of the implemented PDE, coherence enhancing diffusion. We also presented a 3-dimensional implementation as a proof of feasibility.

The construction of the stencils of the AD-LBR is both original and non-trivial. The computational load for this aspect of the algorithm is fortunately not dominant, thanks to the use of a tool from discrete geometry: lattice basis reduction. The AD-LBR also allows to use larger time steps than most of its counterparts, in explicit discretizations of parabolic equations.

AD-LBR trivially extends to vector valued and matrix valued images, by applying it on each image component independently. (In other words, the coupling between image components lies in the construction of the common

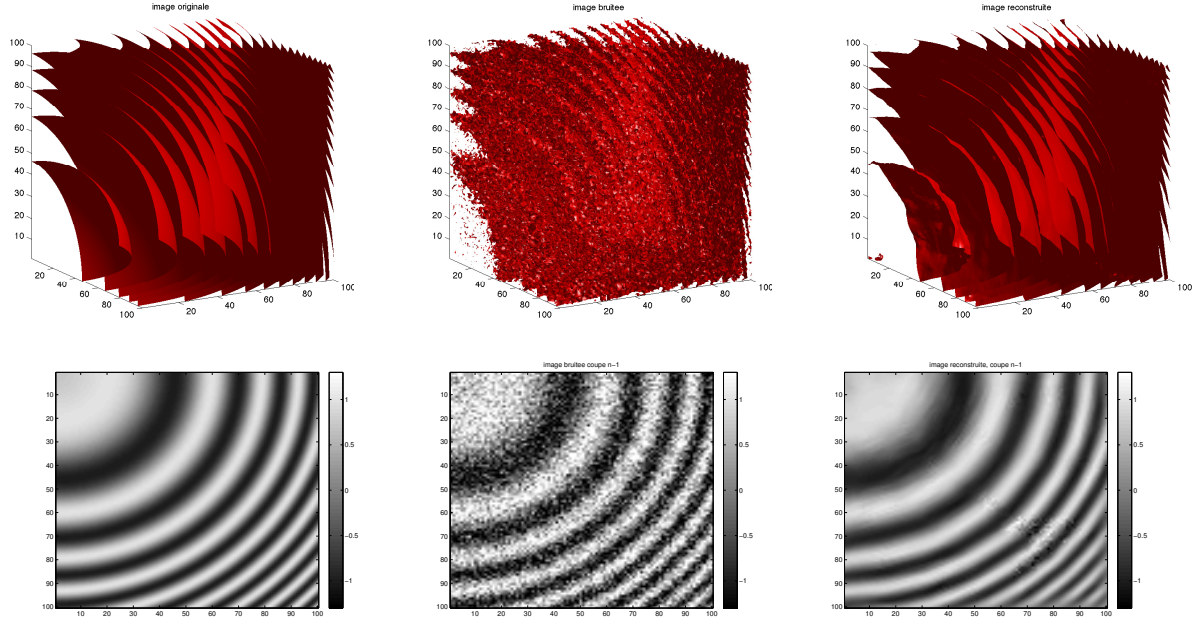


Figure 14: Levelset 0 (top) and slice (bottom) of a 3D image. Original (left), noisy (center), and reconstructed (right) images. Slice in the plane $z = 0.7$, with values clipped to the range $[-1.3, 1.3]$.

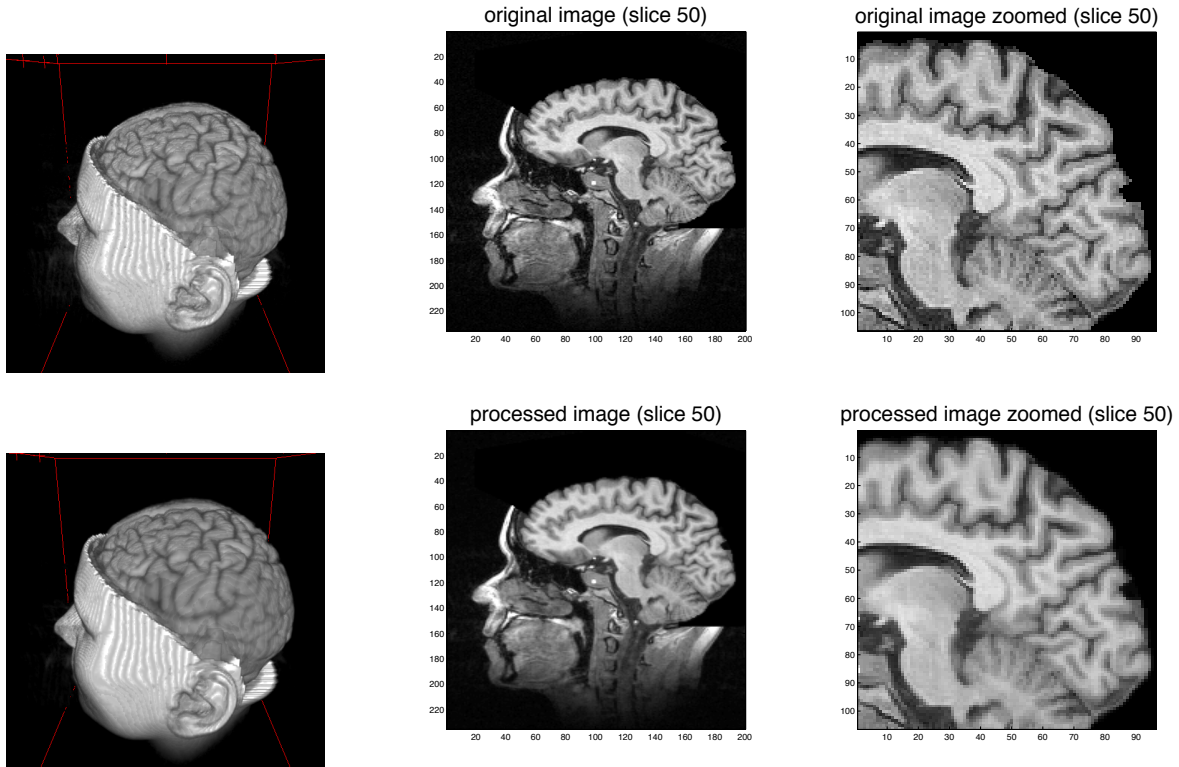


Figure 15: Top: MRI data. Bottom: Data processed via 3D Edge Enhancing diffusion, using AD-LBR. Left: 3D rendering of the original volume data and the processed volume, the 3D rendering was obtained using ImageJ 3D viewer [21] (the effect of anisotropic diffusion is not much visible in this first representation). Center: slice of the original and processed volume. Right: details of the original and processed slices.

diffusion tensor \mathbf{D} , which AD-LBR regards as user input.) Future work will be devoted to the application of AD-LBR to the regularization of diffusion tensor fields, arising for instance from diffusion MRI, for which we expect it to be particularly appropriate: thanks to the scheme non-negativity, positive-definiteness is naturally preserved.

References

- [1] L. Alvarez, F. Guichard, P.-L. Lions, J.-M. Morel, *Axioms and fundamental equations of Image processing*, Arch. Rational Mech. Anal., vol. 123, 199–257 (1993)
- [2] D. Barash, T. Schlick, M. Israeli, and R. Kimmel, *Multiplicative operator splittings in non-linear diffusion: from spatial splitting to multiplicative timesteps*, Journal of Mathematical Imaging and Vision, 19:33–48, (2003).
- [3] J.-B. Bost and K. Künnemann, *Hermitian vector bundles and extension groups on arithmetic schemes. I. Geometry of numbers*, 2010
- [4] J. H. Conway, N. J. A. Sloane, *Low-dimensional lattices. VI. Voronoi reduction of three-dimensional lattices.*, Proceedings of the Royal Society of London. Series A: Mathematical and Physical Sciences 436.1896 (1992): 55–68.
- [5] G.-H. Cottet, L. Germain, *Image processing through reaction combined with nonlinear diffusion*, Math. Comp., vol. 61, 659–673 (1993).
- [6] L. Dascal, A. Ditzkowski, and N. Sochen, *On the Discrete Maximum Principle for the Beltrami Color Flow*, J. Math. Imaging Vision, vol. 29, 63–77 (2007)
- [7] P. Degond, A. Lozinski, J. Narski, C. Negulescu, *An Asymptotic-Preserving method for highly anisotropic elliptic equations based on a micro-macro decomposition*, J. Comput. Phys., vol. 231(7), 2724–2740 (2012)
- [8] S. Grewenig, J. Weickert, and A. Bruhn, *From box filtering to fast explicit diffusion*, Pattern Recognition 533–542, (2010).
- [9] W. Huang, *Discrete maximum principle and a Delaunay-type mesh condition for linear finite element approximations of two-dimensional anisotropic diffusion problems*, arXiv preprint arXiv:1008.0562, (2010)
- [10] F. Labelle, J. R. Shewchuk, *Anisotropic Voronoi Diagrams and Guaranteed-Quality Anisotropic Mesh Generation*, Proceedings of the Nineteenth Annual Symposium on Computational Geometry, 191–200 (2003)
- [11] J. L. Lagrange, *Recherches d’arithmétique*, Nouveaux Mémoires de l’Académie de Berlin, (1773)
- [12] A. K. Lenstra, H. W. Lenstra, and L. Lovász, *Factoring polynomials with rational coefficients*, Mathematische Annalen 261, 513–534, (1982)
- [13] A. M. Mendrik, E. J. Vonken, A. Rutten, M. A. Viergever, and B. van Ginneken, *Noise reduction in computed tomography scans using 3-d anisotropic hybrid diffusion with continuous switch*, Medical Imaging, IEEE Transactions on, 28(10), 1585–1594. (2009)
- [14] J.-M. Mirebeau, *Anisotropic Fast Marching on Cartesian Grids, using Lattice Basis Reduction*, preprint, 2012.
- [15] A. Mitchell and D. Griffiths *The Finite Difference Method in Partial Differential Equations*. Chichester: Wiley (1980).
- [16] P. Mrázek and M. Navara, *Consistent positive directional splitting of anisotropic diffusion*, Proc. Sixth Computer Vision Winter Workshop, (2001)
- [17] P. Q. Nguyen, and D. Stehlé, *Low-dimensional lattice basis reduction revisited*, ACM Transactions on Algorithms, Article 46 (2009).
- [18] P. Q. Nguyen and J. Stern, *The two faces of lattices in cryptology*, In Proceedings of the 2001 Cryptography and Lattices Conference (CALC’01). Lecture Notes in Computer Science, vol. 2146. Springer-Verlag, 146–180, (2001)
- [19] S. Osher, L. Rudin *Feature-oriented image enhancement using shock filters*, SIAM J. Numer. Anal., vol. 27, 919–940 (1990)
- [20] P. Perona and J. Malik, *Scale-Space and Edge Detection Using Anisotropic Diffusion*, IEEE Trans. Patt. Anal. Mach. Int., vol. 12, 629–639 (1990)
- [21] B. Schmid, J. Schindelin, A. Cardona, M. Longair, M. Heisenberg, *A high-level 3D visualization API for Java and ImageJ*, BMC Bioinformatics, 11:274 (2010)
- [22] E. Selling. *über die binären und ternären quadratischen formen*, J. reine angew. Math., 77:143–229, (1874).
- [23] I. Semaev, *A 3-dimensional lattice reduction algorithm*, In Proceedings of the 2001 Cryptography and Lattices Conference (CALC’01). Lecture Notes in Computer Science, vol. 2146. Springer-Verlag, 181–193, (2001)
- [24] J. Weickert, *Theoretical foundations of anisotropic diffusion in image processing*, Computing, SUPPL. 11, 221–236 (1996)
- [25] J. Weickert, *Anisotropic Diffusion in Image Processing*, Teubner, Stuttgart (1998)

- [26] J. Weickert, B. Romeny, and M. Viergever, *Efficient and reliable schemes for nonlinear diffusion filtering*, IEEE Trans. Image Proc., vol. 7, 398–410 (1998)
- [27] J. Weickert, and H. Schar, *A scheme for coherence-enhancing diffusion filtering with optimized rotation invariance*, J. Visual Comm. Image Rep., Vol. 13, 103–118, (2002).
- [28] M. Welk and G. Steidl and J. Weickert, *Locally analytic schemes: A link between diffusion filtering and wavelet shrinkage*, Applied and Computational Harmonic Analysis, 24, pp 195–224, (2008)

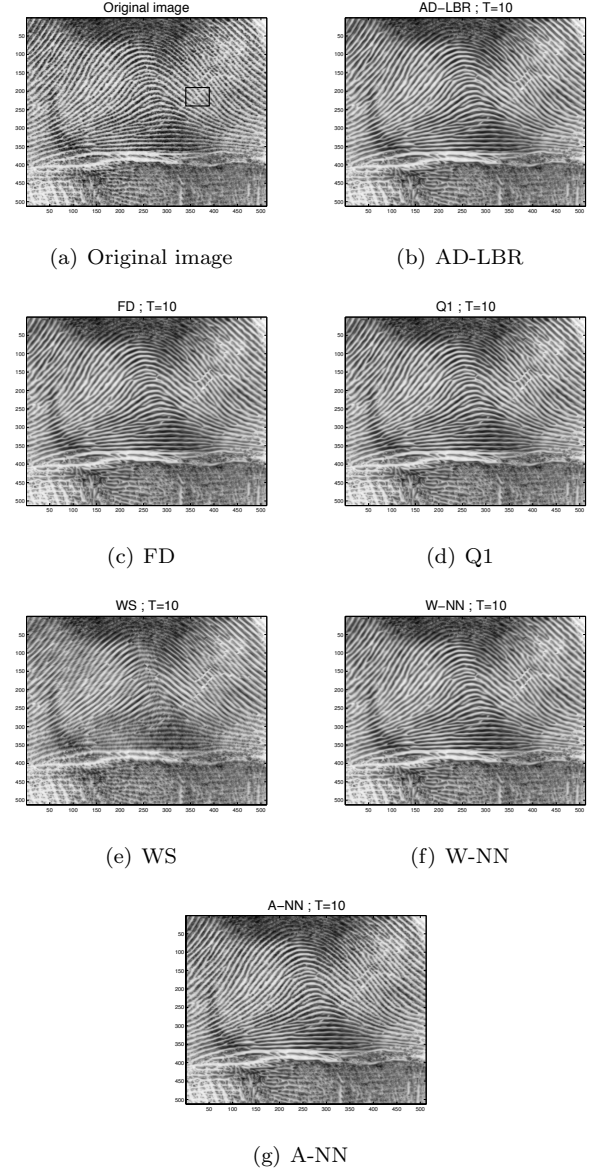


Figure 11: From top to bottom and from left to right: Original image (with two regions highlighted); diffused image using AD-LBR; FD; Q1; WS; W-NN; A-NN. Here $T = 10$.

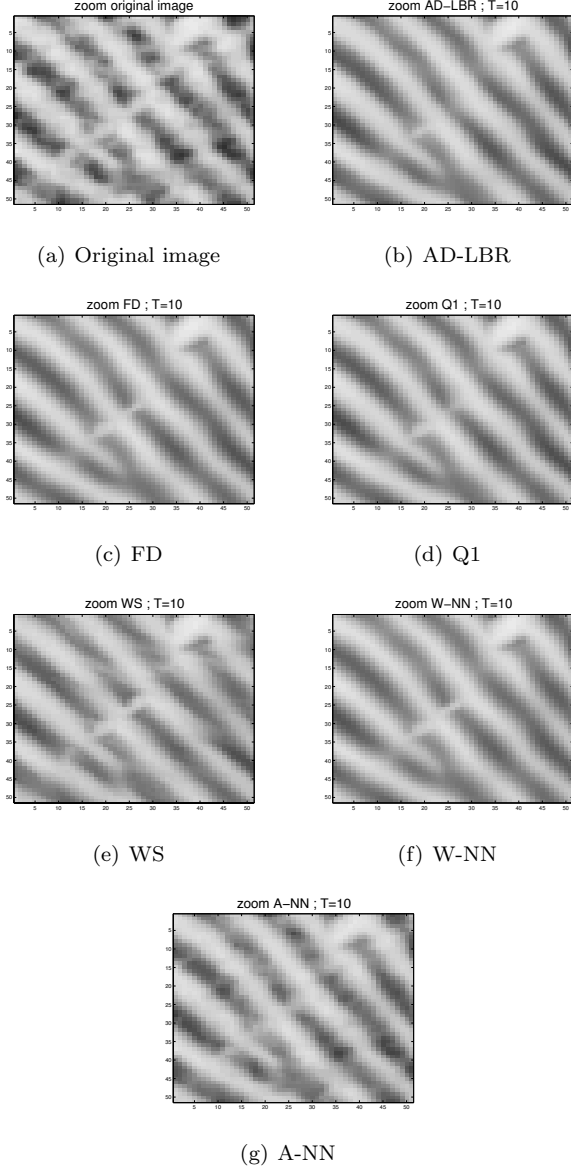


Figure 12: Detail of the region on the right. From top to bottom and from left to right: original image; diffused image using AD-LBR; FD; Q1; WS; W-NN; A-NN.

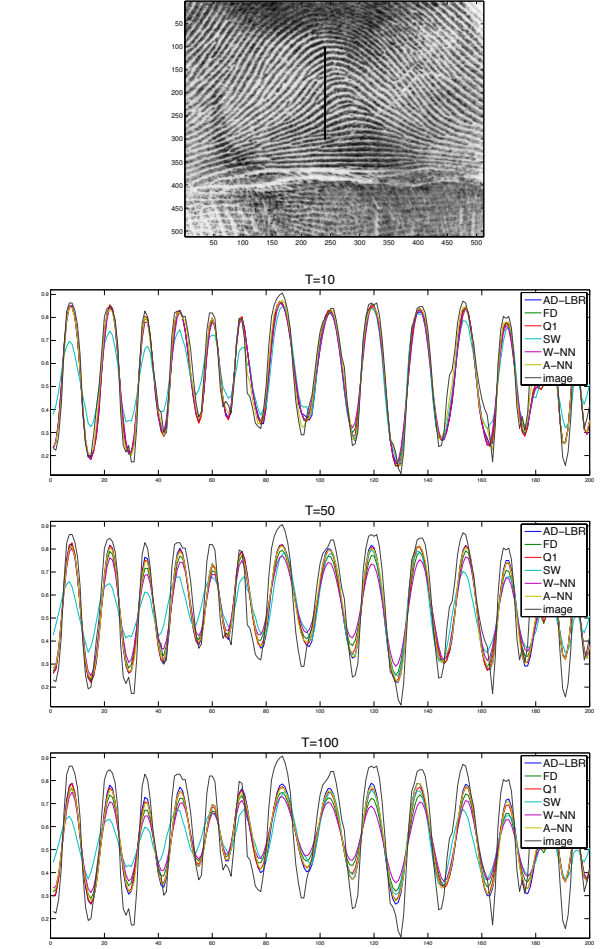


Figure 13: Evolution under CED of a section of the fingerprint image. The ridges in the evolved image are higher, and the valleys are deeper, with AD-LBR than with the other schemes. This illustrates the fact that AD-LBR, respecting the continuous PDE, diffuses more along the structure and less in the orthogonal direction. From top to bottom: location of the section of the image; section at $T = 10$; section at $T = 50$; section at $T = 100$.